Inside MySQL: InnoDB Storage Engine, Second Edition

# MySQL技术内幕
## InnoDB存储引擎

### 第2版

姜承尧◎著

华章经管图书

# MySQL技术内幕

——InnoDB存储引擎 第2版

姜承尧　著

ISBN：978-7-111-42206-8

# 目　录

# 🔲🔲🔲

It's fair to say that MySQL is the most popular open source database.It has a very large installed base and number of users.Let's see what are the reasons MySQL is so popular,where it stands currently,and maybe touch on some of its future(although predicting the future is rarely successful).

Looking at the customer area of MySQL,which includes Facebook,Flickr,Adobe(in Creative Suite 3),Drupal,Digg,LinkedIn,Wikipedia,eBay,YouTube,Google AdSense(source http://mysql.com/customers/and public resources),it's obvious that MySQL is everywhere.When you log in to your popular forum(powered by Bulleting)or blog(powered by WordPress),most likely it has MySQL as its backend database.Traditionally,two MySQL's characteristics,simplicity of use and performance,were what allowed it to gain such popularity.In addition to that,availability on a very wide range of platforms(including Windows)and built-in replication,which provides an easy scale-out solution for read-only clients,gave more user attractions and production deployments.There is simple evidence of MySQL's simplicity:In 15 minutes or less,you really can get installed,have a working database,and start running queries and store data.From its early stages MySQL had a good interface to most popular languages for Web development-PHP and Perl,and also Java and ODBC connectors.

There are two best known storage engines in MySQL:MyISAM and InnoDB(I don't cover NDB cluster here;it's a totally different story).MyISAM comes as the default storage engine and historically it is the oldest,but InnoDB is ACID compliant and provides transactions,row-

level locking,MVCC,automatic recovery and data corruption detection.This makes it the storage engine you want to choose for your application.Also,there is the third-party transaction storage engine PBXT,with characteristics similar to InnoDB,which is included in the MariaDB distribution.

MySQL's simplicity has its own drawback.Just as it is very easy to start working with it,it is very easy to start getting into trouble with it.As soon as your website or forum gets popular,you may figure out that the database is a bottleneck,and that you need special skills and tools to fix it.

The author of this book is a MySQL expert,especially in InnoDB storage engineB.Hence,I highly recommend this book to new users of InnoDB as well as uers who already have well-tuned InnoDB-based applications but need to get internal out of them.

Vadim Tkachenko

知名的MySQL性能优化公司Percona公司CTO

知名MySQL性能博客MySQLPerformanceBlog.com作者

《高性能MySQL》（第2版）作者之一

# □□

## □□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□

□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Think Different□Think Different□20□□90□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Think Differently□□Think Different□Different□□□□□□□□□□□□□□□□□

□□□DBA□□□□□□□□□□□□□"□□"□□□□□□□□"□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Why□□What□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□ □□□□□□□□□"□□"□□□□□□□□□□

□ □□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□ □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□ □□□□□□□□□□□□□□□□□□□□□

可能对于MySQL的初学者有这样的疑问，学习MySQL就是学习SQL语言，就是学习各种各种的数据库操作命令，没有必要了解MySQL数据库底层的相关知识。有这种疑问的读者，大多是还没有遇到各种各种实际问题的困扰。

MySQL数据库独有的插件式存储引擎使得想要成为一名专业的数据库开发和运维人员至关重要。不了解MyISAM和底层实现，不了解InnoDB存储引擎的底层实现，一名DBA可能根本无法判断各种故障发生的原因。一名专业的MySQL DBA，要想更好地运维管理好数据库，深入理解MySQL数据库的底层实现，了解InnoDB存储引擎的各种特性是一条不可忽略也无法跨越的道路。

在明白了学习底层的Why（为What）之后，接下来的问题就是对于InnoDB存储引擎各种底层实现原理的学习要达到一个什么样的程度。笔者认为，对于InnoDB存储引擎底层实现的学习，要达到能够灵活运用的程度。

# 第1版序和第2版序对比

在撰写第2版的过程中，我始终问自己一个问题：是否有足够的内容可以更新？MySQL 5.6版本对于InnoDB存储引擎做了大量的更新，希望将这些内容分享给广大的读者朋友。同时，第1版内容中也存在一些不足，需要加以修正和完善。

相对于第1版，本次第2版的内容更新或修正超过50%，希望读者朋友能够从新版的内容中获得更多的知识。

□全书大约四分之一的内容更新到了MySQL 5.6版本的InnoDB存储引擎。因为MySQL 5.6版本对于InnoDB存储引擎做了大量的更新，InnoDB存储引擎的版本也从之前的版本提升到了新的版本。所以新版中所有的内容都尽可能地更新到了最新的版本，希望读者朋友能够了解到最新版本InnoDB存储引擎的一些特性以及内部的实现原理。

□重写了第七章事务中InnoDB存储引擎redo日志和undo日志部分的内容，修正了第1版中关于InnoDB存储引擎事务中undo部分的一些错误，并结合InnoSQL对于事务部分进行的优化，希望读者对于undo有着更为深刻的认识，从而让DBA对于InnoDB存储引擎有着更深入的了解。

□对第6章锁部分内容进行了大量的修改，重写了InnoDB存储引擎关于next-key locking锁算法部分的内容，修正了之前一些不太正确的部分。同时讲解了InnoDB存储引擎关于锁实现的一些内部原理。

□重写了关于索引中InnoDB存储引擎insert buffer部分的内容，因为新版本对于这部分内容的实现原理也进行了修改，原先的insert buffer已经升级为了——change buffer的技术。

# 内容简介

本书是数据库内核开发者呈现给广大数据库从业者的一份厚礼。全书以MySQL数据库SQL语句的执行流程为主线，以MySQL InnoDB存储引擎的实现为核心，深入浅出地介绍了MySQL数据库的各个方面。无论是MySQL DBA，还是数据库内核开发者，阅读SQL语句的执行流程，都能从本书中汲取知识，提升对MySQL的理解。从Replication原理，到数据库锁机制，再到事务处理，InnoDB存储引擎的内部实现，本书都做了深入的剖析，带领读者一步步走进数据库的内核世界。

本书的最大特点是，以InnoDB存储引擎的实现为核心，结合源码，深入浅出地剖析数据库内部的实现原理。读者在阅读本书的过程中，不仅能够掌握数据库的使用方法，更能够理解MySQL数据库内部的实现机理，做到知其然，更知其"所以然"，从而在实际工作中游刃有余。

本书适合以下读者阅读，具体来说，包括：

□ 精通SQL，

□ 深入理解MySQL原理，

□ 熟悉一门编程语言，如C、C++、Python、Java。

□ 对数据库内核感兴趣，希望深入理解InnoDB存储引擎的实现原理，提升自身技术水平的数据库从业者。

# 如何阅读本书

本书的前面10章基本来自我自己的"老本行"——数据库内核开发。这些内容相对独立，可以单独阅读。第10章是对InnoDB存储引擎文件的详细介绍，可以将它作为前3章的附加章节。也许你在读前几章时会对表产生疑问，希望第10章能够解答你心中的疑问。

对于其他的读者，也可以将本书作为一本参考手册。如果你对InnoDB存储引擎的某部分内容感兴趣，可以根据相应的章节直接查阅。此外，通过API你也能了解InnoDB存储引擎内部的很多机制，这对排查问题来说可能会有帮助。

对于书中的示例，可以到华章网站的本书页面（www.hzbook.com）进行下载。如果对本书有任何疑问或勘误，可以给我发送邮件。

# 写在最后

由于本人在InnoDB存储引擎方面的知识有限，因此书中难免会有描述不准确或者错误的地方，希望广大读者能够指正，大家相互讨论，相互进步。我的电子邮件地址是jiangchengyao@gmail.com，新浪微博账号是@insidemysql。如果你对书中的任何问题有任何的疑问，都可以和我联系，我将非常感谢!

# 序言

本书及其姉妹书积累了许多作者的经验。翻译过程中，作者之一Pecona公司的CEO Peter Zaitsev、CTO Vadim Tkachenko等也给予了大量帮助。本书主要讲解InnoDB相关内容，大部分是基于作者多年的工作经验而写就，对于InnoDB存储引擎的讲解非常深入。

本书内容丰富，讲解透彻，不仅适合数据库开发人员阅读，也适合数据库管理人员以及对数据库底层原理感兴趣的读者阅读。

在本书的翻译过程中，得到了许多朋友的帮助，在此一并表示感谢。由于译者水平有限，书中难免存在不足之处，敬请读者批评指正。

最后，感谢家人在翻译过程中对我的理解和支持，感谢所有帮助过我的朋友，也感谢广大读者对本书的关注和厚爱。

译 者

# 第1章 MySQL用户与权限管理

MySQL数据库是一个多用户数据库，并且支持多平台操作，包括Linux、Solaris、FreeBSD、Mac、Windows等操作系统平台都可以很好地运行。本章主要介绍MySQL数据库用户与权限的管理操作，通过本章的学习，读者应掌握对MySQL数据库用户的创建与权限的管理方法。

# 1.1 □□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□"□□□"（database）□"□□"（instance）□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□frm、MYD、MYI、ibd□□□□□□□□□□NDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

MySQL□□□□□□□□□□□□□□□□□□□□□□□SQL Server□□□□□□□□Oracle□□□□□□□□□□□Oracle、Windows□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□

□Linux□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□ps□□MySQL□□□□□□□□□□□□□□

---

```
[root@xen-server bin]#./mysqld_safe□

[root@xen-server bin]#ps-ef|grep mysqld

root 3441 3258 0 10:23 pts/3 00:00:00/bin/sh./mysqld_safe

mysql 3578 3441 0 10:23 pts/3 00:00:00

/usr/local/mysql/libexec/mysqld--basedir=/usr/local/mysql

--datadir=/usr/local/mysql/var--user=mysql

--log-error=/usr/local/mysql/var/xen-server.err

--pid-file=/usr/local/mysql/var/xen-server.pid

--socket=/tmp/mysql.sock--port=3306

root 3616 3258 0 10:27 pts/3 00:00:00 grep mysqld
```

□□□□□□□3578□□□□□□□□□MySQL□□□□□□□□□□□□□mysqld_safe□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Oracle□□□□□□□spfile□□□□□□□□□□Oracle□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□

```
[root@xen-server bin]#mysql--help|grep my.cnf

order of preference,my.cnf,$MYSQL_TCP_PORT,

/etc/my.cnf/etc/mysql/my.cnf/usr/local/mysql/etc/my.cnf□/.my.cnf
```

□□□□□□MySQL□□□□□□
□/etc/my.cnf→/etc/mysql/my.cnf→/usr/local/mysql/etc/my.cnf→□/.my.cnf□□□□□□□□□□□□□□□□□□□□□□□□"□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□"□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Linux□□□□□□□□□□□□□□□/etc/my.cnf□□□□Windows□□□□□□□□□□□□□□□□□□.cnf□□□□□□.ini□□□□Windows□□□□□□□□mysql--help□□□□□□□□□□□□□□□□□

```
Default options are read from the following files in the given order:

C:\Windows\my.ini C:\Windows\my.cnf C:\my.ini C:\my.cnf C:\Program Files\MySQL\M

\MySQL Server 5.1\my.cnf
```

□□□□□□□□□□□□datadir□□□□□□□□□□□□□□□□□□□□□□□□Linux□□□□□□□□□datadir□/usr/local/mysql/data□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
mysql□SHOW VARIABLES LIKE'datadir'\G;

***************************1.row***************************
```

Variable_name:datadir

Value:/usr/local/mysql/data/

1 row in set(0.00 sec)1 row in set(0.00 sec)

mysql〔system ls-lh/usr/local/mysql/data

total 32K

drwxr-xr-x 2 root mysql 4.0K Aug 6 16:23 bin

drwxr-xr-x 2 root mysql 4.0K Aug 6 16:23 docs

drwxr-xr-x 3 root mysql 4.0K Aug 6 16:04 include

drwxr-xr-x 3 root mysql 4.0K Aug 6 16:04 lib

drwxr-xr-x 2 root mysql 4.0K Aug 6 16:23 libexec

drwxr-xr-x 10 root mysql 4.0K Aug 6 16:23 mysql-test

drwxr-xr-x 5 root mysql 4.0K Aug 6 16:04 share

drwxr-xr-x 5 root mysql 4.0K Aug 6 16:23 sql-bench

lrwxrwxrwx 1 root mysql 16 Aug 6 16:05 data-〔/opt/mysql_data/

# □□□□□□□□□□data□□□□□□□□□□□□□□□□□/opt/mysql_data□□□□□□□□□□□□□□/opt/mysql_data□□□□□□□□□□□mysql□□□□□□□□□□□□MySQL□□□□□□□mysql:mysql□□

## 1.2 MySQL之选择

每次在我的微信公众号或者学习群里面有初学者问我，我想学习数据库，我想成为一名DBA，那么我应该学习什么数据库。每当有人问到MySQL的时候，总有些人会觉得学习MySQL没有前途，不像SQL Server、Oracle、DB2那么高大上，问的最多的几个问题就是：

□为什么MySQL有那么多分支？

□MySQL都被收购了，还有前途吗？

□国内有超过1000家的MySQL使用厂商吗？

……

所以MySQL真的像他们说的那样没有前途吗？想必现在大家的心里都会有一个疑问，下面我就带着这个疑问为大家揭晓问题的答案。

如果各位是一名初学者，相信大家在学习一门技术或者一门课程之前都会考虑为什么要学习MySQL，学习了MySQL有什么用？等诸如此类的问题，不像Oracle、Microsoft SQL Server那样无须大家质疑它们的强大，反而MySQL给大家的第一印象就是，它是一个免费开源的数据库，相比其他商业数据库而言显得很渺小，甚至有的人会觉得学习它根本就没有什么太大的价值，其实不然，我只能说会有这样的想法的读者一定不了解什么是开源数据库，更不了解它的重要性和在实际生产环境中的应用，下面我将带领大家去认识它。

其实数据库大体可以分为两类，一类是我们付费使用的商业数据库，另一类是我们可以免费使用的开源数据库，它们都支持SQL，如SELECT、INSERT、UPDATE、DELETE等，也就是说它们都是关系型数据库，本质上基本没有太大的区别，如果你细心可能会发现它们中间的一些差异，如Oracle、SQL Server、MySQL在分页和获取系统时间等函数的使用上存在着一些差异，但是这丝毫不会影响我们的使用。

接下来我将通过一系列的数据和图表的方式来给大家展示MySQL的强大和重要性，如图1-1所示，展示的是MySQL近些年的发展。

Connectors
Native C API, JDBC, ODBC, NET, PHP, Perl, Python, Ruby, Cobol

**MySQL Server**

Management Service & Utillties

Backup & Recovery, Security, Replication, Cluster, Administration Configuration, Migration & Metadata

Connection Pool
Authentication, Thead Reuse, Connection Limits, Check Memory, Caches

SQL Interface
DML, DDL, Stored Procedures Views, Triggers, etc.

Parser
Quary Translation Object Privilege

Optimizer
Access Paths, Statistics

Cathes & Buffers
Global and Engine Specific Caches & Buffers

Pluggable Storage Engines
Memory, Index & Storage Management

MyISAM   InnoDB   NDB   Archive Federated Memory   Merge   Partner Community Custom

File system
NTFS, ufs, ext2/3
NFS, SAN, NAS

Files & Logs
Redo, Undo, Data, Index,
Binary, Error, Query and Slow

## 图 1-1　MySQL逻辑结构

如图1-1所示描述了MySQL的逻辑结构，主要包括：

□连接管理器

□线程管理及线程池

□SQL分析处理

□查询缓存管理器

□存储引擎

□线程及Cache管理器

□用户认证管理器

□内存管理

如图1-1所示描述了MySQL的逻辑结构，它是由多个层次和多个组件所构成的。在这些组件中，MySQL的查询缓存管理器、连接管理器、线程管理器、用户认证管理器是属于逻辑结构上层的组件。SQL分析处理、查询缓存管理器、存储引擎等组件属于逻辑结构中的核心组件。

为了让读者能更清楚地了解每一个组件在逻辑结构中的位置，图1-1中MySQL的逻辑结构清晰地描述了MySQL逻辑结构中各组件的关系。

# 1.3 MySQL□□□□

□□1.2□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□

□□MySQL□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□eBay□□□□□Igor Chernyshev□MySQL Memory□□□□□□□□□□http://code.google.com/p/mysql-heap-dynamic-rows/□□□□□□eBay□Personalization Platform□□□□□□□□□Google□Facebook□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□SSD□□□□□□□ [1]□□□□□□□SSD□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□16□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□MySQL□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□Oracle□□□□□□□□□□□□□□□□□□□MySQL□□□□OLTP□Online Transaction Processing□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 1.3.1 InnoDB□□□□

InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLTP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Oracle□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□5.5.8□□□□□□InnoDB□□□□□□□□□□□□□□□□

InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□MySQL 4.1□□□4.1□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□ibd□□□□□□□□□InnoDB□□□□□□□□□□□□□row disk□□□□□□□□□□□□□□

InnoDB□□□□□□□□□□□□□□MVCC□□□□□□□□□□□□□□□□□□SQL□□□□4□□□□□□□□□□REPEATABLE□□□□□□□□□□□□□□□next-key locking□□□□□□□□□□

和phantom现象带来的问题。此外，InnoDB存储引擎还提供了插入缓冲（insert buffer）、二次写（double write）、自适应哈希索引（adaptive hash index）、预读（read ahead）等高性能和高可用的功能。

对于表中数据的存储，InnoDB存储引擎采用了聚集（clustered）的方式，因此每张表的存储都是按主键的顺序进行存放。如果没有显式地在表定义时指定主键，InnoDB存储引擎会为每一行生成一个6字节的ROWID，并以此作为主键。

InnoDB存储引擎是MySQL数据库最为常用的一种引擎，而Facebook、Google、Yahoo等公司的成功应用已经证明了InnoDB存储引擎具备的高可用性、高性能以及高可扩展性。

[1]____参     见     网     址     http://code.google.com/p/david-mysql-tools/wiki/innodb_secondary_buffer_pool

# 1.3.2 MyISAM□□□□

MyISAM□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLAP□□□□□□□□□
MySQL 5.5.8□□□□□MyISAM□□□□□□□□□□□□□□□□□□□□□Windows□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MyISAM□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ETL□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□MyISAM□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□cache□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

MyISAM□□□□□□□MYD□MYI□□□□MYD□□□□□□□□□□MYI□□□□□□□□□□□□□□
□□□□□myisampack□□□□□□□□□□□□□□□□□myisampack□□□□□□□□□□
□Huffman□□□□□□□□□□□□□□□□□□□□□□myisampack□□□□□□□□□□□□□□□□□
□□□□□□□myisampack□□□□□□□□□□□

□MySQL 5.0□□□□□□□MyISAM□□□□□□□□□□□4GB□□□□□□□□□□□□4GB□
MyISAM□□□□□□□□□□□□MAX_ROWS□AVG_ROW_LENGTH□□□□□□
MySQL 5.0□□□□□□□MyISAM□□□□□256TB□□□□□□□□□□□□□□□□□□□□□□□

□□ □□MyISAM□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□LRU□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL 5.1.23□□□□
□□□□□□□□32□□□64□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□4GB□□□□□□□□□□□
64□□□□□□□□□□□4GB□□□□□□□□□□

### 1.3.3 NDB存储引擎

2003年，MySQL AB公司从Sony Ericsson公司收购了NDB存储引擎。图1-1显示了NDB存储引擎的架构，可以看出，与Oracle的RAC不同，其并不是Oracle RAC share everything架构，而是share nothing的集群架构，因此能提供更高级别的高可用性。NDB的特点是数据全部放在内存中（从MySQL 5.1版本开始，可以将非索引数据放在磁盘上），因此主键查找（primary key lookups）的速度极快，并且能够在线添加NDB数据存储节点（Data Node）以便线性地提高数据库性能。由此可见，NDB是高可用、高性能、高可扩展性的数据库集群系统，其面向的也是OLTP的数据库应用类型。

虽然NDB看起来是个完美的集群系统，但其有一个问题值得注意：连接操作（JOIN）是在MySQL数据库层完成的，而不是在存储引擎层完成的。这意味着，复杂的连接操作需要巨大的网络开销，因此查询速度很慢。如果解决了这个问题，NDB存储引擎的市场应该是非常巨大的。

注意　MySQL NDB Cluster是一个单独的软件包，用户可以在其官网上下载NDB Cluster版本（Carrier Grade Edition）。有兴趣的读者可以访问http://dev.mysql.com/downloads/cluster/index.html，以了解更多关于NDB Cluster的相关信息。

## 1.3.4 Memory存储引擎

Memory存储引擎（之前称HEAP存储引擎）将表中的数据存放在内存中，如果数据库重启或发生崩溃，表中的数据都将消失。它非常适合用于存储临时数据的临时表。Memory存储引擎默认使用哈希索引，而不是我们熟悉的B+树索引。

虽然Memory存储引擎速度非常快，但在使用上还是有一定的限制。比如，只支持固定长度的行，TEXT和BLOB等字段类型不被支持。因此哪怕用户将varchar类型定义，实际使用的还是char类型，这对内存的浪费显得尤为明显（这个问题，现eBay的架构师Igor Chernyshev曾给出过patch来解决）。

此外还要注意的是，MySQL数据库使用Memory存储引擎作为临时表来存放查询的中间结果集（intermediate result）。如果中间结果集大于Memory存储引擎表的容量设置，又或者中间结果集含有TEXT或BLOB列类型字段，则MySQL数据库会把其转换到MyISAM存储引擎表而存放到磁盘中。之前提到MyISAM不缓存数据文件，因此这时产生的临时表的性能对于查询会有损失。

# 1.3.5　Archive□□□

Archive□□□□□□□INSERT□SELECT□□□□MySQL 5.1□□□□□□□□Archive□□□□□□zlib□□□□□□□□row□□□□□□□□□□□□□□□□□1:10□□□□□□□□□□Archive□□□□□□□□□□□□□□□□□□□□□□□□□□Archive□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 1.3.6 Federated存储引擎

Federated存储引擎是访问其他MySQL服务器的一个代理，尽管它看起来提供了一种很好的跨服务器的灵活性，但也经常带来问题，因此默认是禁用的。

## 1.3.7　Maria□□□□

Maria□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MyISAM□□□□□□□□□□□MySQL□□□□□□□□□□Maria□□□□□□□□□□MySQL□□□□□□□□Michael Widenius□□□□□□□□□□□□MyISAM□□□□□□□Maria□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MVCC□□□□□□□□□□□□□□□□□□□□□□□□□□□□BLOB□□□□□□□□□□□□

# 1.3.8 □□□□□□□

□□□□□□□7□□□□□□□MySQL□□□□□□□□□□□□□□□□□□Merge□CSV□Sphinx□Infobright□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□1.2□□□□□□□□□□

□□□□MySQL□□□□□□□□□□□□□□□□MySQL□□□MyISAM□InnoDB□1.2□□□□□Sphinx□□□□□□□□□□□□□□

□MySQL□□□□□□□□□□□□□□□□□□□□MySQL□MyISAM□□□□□□□□□□□□□□InnoDB□□□□"□"□□□□□□□□□□□□□□□□ETL□□□□□□MyISAM□□□□□□□□□□OLTP□□□□InnoDB□□□□□□□□□□□□□

□□□□□□□□□1000□□MySQL□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□Mytrix□Inc.□InnoDB□□□□□□1 TB□□□□□□□□□□□□□□□InnoDB□□□□□□□□/□□□□□□□800□/□□

## 1.4 本章小结与下章概要

经过1.3节的讲解，相信读者已经对MySQL的基础架构有了大致的了解，这对于后面章节的学习是非常有帮助的。在1-2节中，MySQL的事务底层没有讲，因为MySQL的事务底层还没有涉及到，在后面的章节中会详细讲解，其中就包括MVCC的实现原理等，欢迎大家持续关注。

| Feature | MyISAM | BDB | Memory | InnoDB | Archive | NDB |
|---|---|---|---|---|---|---|
| Storage Limits | No | No | Yes | 64TB | No | Yes |
| Transactions (commit, rollback, etc.) | | ✔ | | ✔ | | |
| Locking granularity | Table | Page | Table | Row | Row | Row |
| MVCC/Snapshot Read | | | | ✔ | ✔ | ✔ |
| Geospatial support | ✔ | | | | | |
| B-Tree indexes | ✔ | ✔ | ✔ | ✔ | | ✔ |
| Hash indexes | | | ✔ | ✔ | | ✔ |
| Full text search index | ✔ | | | | | |
| Clustered index | | | | ✔ | | |
| Data Caches | | | ✔ | ✔ | | ✔ |
| Index Caches | ✔ | | ✔ | ✔ | | ✔ |
| Compressed data | ✔ | | | | ✔ | |
| Encrypted data (via function) | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |
| Storage cost (space used) | Low | Low | N/A | High | Very Low | Low |
| Memory cost | Low | Low | Medium | High | Low | High |
| Bulk Insert Speed | High | High | High | Low | Very High | High |
| Cluster database support | | | | | | ✔ |
| Replication support | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |
| Foreign key support | | | | ✔ | | |
| Backup/Point-in-time recovery | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |
| Query cache support | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |
| Update Statistics for Data Dictionary | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |

# 图 1-2 查看MySQL数据库所支持的引擎

数据库引擎是用于存储、处理和保护数据的核心服务。利用数据库引擎可以控制访问权限并快速处理事务，从而满足企业内大多数需要处理大量数据的应用程序的要求。MySQL存储引擎是数据库底层软件组件，数据库管理系统使用MySQL引擎进行创建、查询、更新和删除数据操作。可以使用MySQL提供的如下命令来查看当前数据库所支持的引擎信息。

可以使用SHOW ENGINES语句来查看系统所支持的MySQL引擎类型，或者也可以通过查询information_schema数据库下ENGINES表进行查看。

---

```
mysql>SHOW ENGINES\G;

***************************1.row***************************

Engine:InnoDB

Support:YES

Comment:Supports transactions,row-level locking,and foreign keys

Transactions:YES

XA:YES

Savepoints:YES

***************************2.row***************************

Engine:MRG_MYISAM

Support:YES

Comment:Collection of identical MyISAM tables

Transactions:NO

XA:NO

Savepoints:NO

***************************3.row***************************

Engine:BLACKHOLE

Support:YES

Comment:/dev/null storage engine(anything you write to it disappears)

Transactions:NO

XA:NO

Savepoints:NO
```

```
***************************4.row***************************

Engine:CSV

Support:YES

Comment:CSV storage engine

Transactions:NO

XA:NO

Savepoints:NO

***************************5.row***************************

Engine:MEMORY

Support:YES

Comment:Hash based,stored in memory,useful for temporary tables

Transactions:NO

XA:NO

Savepoints:NO

***************************6.row***************************

Engine:FEDERATED

Support:NO

Comment:Federated MySQL storage engine

Transactions:NULL

XA:NULL

Savepoints:NULL

***************************7.row***************************

Engine:ARCHIVE

Support:YES

Comment:Archive storage engine

Transactions:NO

XA:NO

Savepoints:NO

***************************8.row***************************

Engine:MyISAM
```

Support:DEFAULT

Comment:Default engine as of MySQL 3.23 with great performance

Transactions:NO

XA:NO

Savepoints:NO

8 rows in set(0.00 sec)

---

可以发现，MySQL数据库有着非常多的存储引擎，用户可以根据具体的应用建立于不同存储引擎之上的表，下表显示了一个例子。

---

mysql□CREATE TABLE mytest Engine=MyISAM

-□AS SELECT*FROM salaries;

Query OK,2844047 rows affected(4.37 sec)

Records:2844047 Duplicates:0 Warnings:0

mysql□ALTER TABLE mytest Engine=InnoDB;

Query OK,2844047 rows affected(15.86 sec)

Records:2844047 Duplicates:0 Warnings:0

mysql□ALTER TABLE mytest Engine=ARCHIVE;

Query OK,2844047 rows affected(16.03 sec)

Records:2844047 Duplicates:0 Warnings:0

---

这个示例显示了从存储引擎为MyISAM的表转换到不同存储引擎的表的过程，其中MyISAM存储引擎下表的大小为40.7MB，InnoDB存储引擎下表的大小为113.6MB，而在Archive存储引擎下表的大小为20.2MB。从这个例子中你还能发现哪些有趣的现象呢？

注意　MySQL数据库区别于其他数据库的最重要的一个特点就是其插件式的表存储引擎。MySQL数据库的核心并不是表存储引擎，而是连接器、查询分析器、优化器、缓存等组件构成的逻辑。用户可以根据具体应用的需求在表上选择不同的存储引擎，这好比在SQL Server中有AdventureWorks示例数据库，在Oracle中有scott用户下的一系列表，而在MySQL数据库中你可以根据自己的需求选择合适的存储引擎。想了解更多有关Oracle、SQL Server等数据库存储引擎的知识，读者可参考各数据库的相关手册，如果你对MySQL数据库的存储引擎感兴趣，请参阅http://dev.mysql.com/doc/。

# 1.5 　连接MySQL

本节将介绍如何连接MySQL数据库。连接数据库是进行数据库操作MySQL服务器的前提条件。只有在成功连接数据库之后，才可以进行后续的一系列操作。连接数据库主要有两种方式，即TCP/IP连接和UNIX套接字连接MySQL数据库服务器。下面分别介绍这两种数据库的连接方法。

# 1.5.1 　TCP/IP

TCP/IP连接是通过MySQL客户端程序与服务器建立的基于某个端口的网络连接。数据库用户使用TCP/IP方式建立连接时，该用户就是使用client程序的主机用户，连接到MySQL数据库server端的某个端口。这种方式通过TCP/IP方式建立连接，可以在Windows主机环境中，建立对Linux主机上的MySQL服务器访问。

---

```
C:\>mysql-h192.168.0.101-u david-p

Enter password:

Welcome to the MySQL monitor.Commands end with;or\g.

Your MySQL connection id is 18358

Server version:5.0.77-log MySQL Community Server(GPL)

Type'help;'or'\h'for help.Type'\c'to clear the current input statement.

mysql>
```

---

上述语句是从Windows主机登录了Host IP为192.168.0.101的MySQL数据库。在TCP/IP方式建立连接之后，可以对MySQL数据库进行后续的各种DDL、DML操作。

通过网络连接数据库时，TCP/IP连接MySQL数据库的MySQL用户，必须经过数据库的身份验证，同时，该主机的IP地址必须在MySQL数据库系统的mysql数据库下面的user表中，授权可用。

---

```
mysql>USE mysql;

Database changed
```

```
mysql]SELECT host,user,password FROM user;

***************************1.row***************************

host:192.168.24.%

user:root

password:*75DBD4FA548120B54FE693006C41AA9A16DE8FBE

***************************2.row***************************

host:nineyou0-43

user:root

password:*75DBD4FA548120B54FE693006C41AA9A16DE8FBE

***************************3.row***************************

host:127.0.0.1

user:root

password:*75DBD4FA548120B54FE693006C41AA9A16DE8FBE

***************************4.row***************************

host:192.168.0.100

user:zlm

password:*DAE0939275CC7CD8E0293812A31735DA9CF0953C

***************************5.row***************************

host:%

user:david

password:

5 rows in set(0.00 sec)
```

□□□□□□□□□□□□MySQL□□david□□□□□□□□IP□□□□□□□□□□□□□□□□□□□
□□□□□□□root□□□□□□□□□□□□□□□□□□

## 1.5.2　命名管道和共享内存

在Windows 2000、Windows XP、Windows 2003、Windows Vista中，如果客户端和服务器在同一台机器上，那么可以使用命名管道来连接，这和Microsoft SQL Server的命名管道非常类似。但是在使用命名管道时MySQL需要在启动时加上--enable-named-pipe选项。在MySQL 4.1以后的版本中，MySQL还提供了共享内存的连接方式，这是通过--shared-memory实现的。如果想使用共享内存的方式，在启动MySQL服务时还必须添加--protocol=memory选项。

# 1.5.3 UNIX□□□□

□Linux□UNIX□□□□□□□□□□UNIX□□□□□□UNIX□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□--socket=/tmp/mysql.sock□□□□□□□□□□□□□□□□□□□□□□□□□□UNIX□□□□□□□□□□□□

---

mysql□SHOW VARIABLES LIKE'socket';

***************************1.row***************************

Variable_name:socket

Value:/tmp/mysql.sock

1 row in set(0.00 sec)

---

## □□□□□UNIX□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

[root@stargazer□]#mysql-udavid-S/tmp/mysql.sock

Welcome to the MySQL monitor.Commands end with;or\g.

Your MySQL connection id is 20333

Server version:5.0.77-log MySQL Community Server(GPL)

Type'help;'or'\h'for help.Type'\c'to clear the buffer.

mysql□

---

# 1.6   小结

本章循序渐进地介绍了数据库的一些基本概念，以及MySQL的一些基础知识。通过本章的学习，读者可以对"数据库"和"数据库系统"有一个初步认识。作为MySQL DBA，对MySQL的源代码有一定的了解是必要的，掌握MySQL的编译和安装，是MySQL数据库管理最基本的要求。本章还简单介绍了MySQL的各种启动方式，以及它们之间的区别与联系。

下一章将介绍各种常见的存储引擎，比较它们之间的区别，以及它们的应用环境，并指导读者如何选择合适的MySQL存储引擎，为以后数据库的设计和应用打下坚实的基础。存储引擎是数据库的核心，对于开发人员来说，也许不需要非常深入地了解它，但是对于数据库管理人员来说，这一部分是必须掌握的知识。

# 第2章　InnoDB存储引擎

InnoDB是事务安全的MySQL存储引擎，设计上采用了类似于Oracle数据库的架构。通常来说，InnoDB是构建高性能OLTP类型应用数据库系统的首选方案。同样地，InnoDB也是目前MySQL数据库官方版本的默认存储引擎。本章将详细地介绍InnoDB存储引擎的体系架构及其各个组成模块的相关工作原理。

# 2.1 InnoDB□□□□□□

InnoDB□□□□□□□Innobase Oy□□[1]□□□□□□□MySQL□□□□□□□□□□□□□□□□□MySQL 5.5□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□Windows□□□□□□□□□□□□□□□□□□□□□□□□□ACID□□□MySQL□□□□□□BDB□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MVCC□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□CPU□

Heikki Tuuri□1964□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□Linux□□□□Linus□□□□□□□□□□□□□□1990□□□□□□□□□□□□□□□□□□□□□□1995□□□Innobase Oy□□□□□CEO□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□Calvin Sun□□□□□□□□□□□□Jimmy Yang□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Sybase□□□□□□□□□□□□□□□□□□□□

InnoDB□□□□□□□□□□□□□□□□□□□□□□□□Google□Yahoo!□Facebook□YouTube□Flickr□□□□□□□□□□□□□□□□□□□□Second Life□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□OLTP□□□□□MySQL InnoDB□□□□□□□□□□□□□□

□MySQL□□□□□□□□□□□□□□□□Internet□□□□□Slashdot.org□□□□InnoDB□□Mytrix□Inc.□InnoDB□□□□□□1 TB□□□□□□□□□□□□□□□InnoDB□□□□□/□□□□□□□□□□□□800□/□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□

InnoDB□□□□□MySQL□□□□□□□□GNU GPL 2□□□□□□□□□MySQL□□□□□□□□□□http://www.mysql.com/about/legal/□□□□□□□□□□□

[1] 2006□□□□□□□□Oracle□□□□□

## 2.2 InnoDB存储引擎的版本

InnoDB存储引擎被包含于所有MySQL数据库的二进制发行版本中。早期其版本随着MySQL数据库的更新而更新。从MySQL 5.1版本时，MySQL数据库允许存储引擎开发商以动态方式加载引擎，这样存储引擎的更新可以不受MySQL数据库版本的限制。所以在MySQL 5.1中，可以支持两个版本的InnoDB，一个是静态编译的InnoDB版本，可将其视为老版本的InnoDB；另一个是动态加载的InnoDB版本，官方称为InnoDB Plugin，可将其视为InnoDB 1.0.x版本。MySQL 5.5版本中又将InnoDB的版本升级到了1.1.x。而在最新的MySQL 5.6中InnoDB的版本也随着升级到1.2.x。表2-1显示了各个版本InnoDB存储引擎的功能。

表 2-1　InnoDB 各版本功能对比

| 版　　本 | 功　　能 |
| --- | --- |
| 老版本 InnoDB | 支持 ACID、行锁设计、MVCC |
| InnoDB 1.0.x | 继承了上述版本所有功能，增加了 compress 和 dynamic 页格式 |
| InnoDB 1.1.x | 继承了上述版本所有功能，增加了 Linux AIO、多回滚段 |
| InnoDB 1.2.x | 继承了上述版本所有功能，增加了全文索引支持、在线索引添加 |

为了方便理解，这里举例说明MySQL数据库的版本和InnoDB MySQL 5.1中支持两个版本的InnoDB Plugin的关系。DBA可以选择将InnoDB Plugin（InnoDB 1.1）替换掉老版本的，而如表2-1所示，老版本的不支持新的页格式compress和dynamic，同时也不支持InnoDB Plugin所支持的Linux Native AIO。因此，可能会有读者朋友会问：InnoDB

Plugin的功能。另外，本书接下来关注的1023页，它是为MySQL 5.5版本中自带的InnoDB Plugin所进行的介绍和分析。

## 2.3　InnoDB□□□□

□□□1□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□2-1□□□□□InnoDB□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□/□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□redo log□□□□□□

……

图 2-1 InnoDB存储引擎体系结构

保证了数据的一致性。同时，缓冲池还被用来协调对磁盘的读写操作。所以，本书之前讨论的各种数据结构（如缓冲池等），都是构成了InnoDB存储引擎的一个整体结构。

# 2.3.1 后台线程

InnoDB存储引擎是多线程的模型，因此其后台有多个不同的后台线程，负责处理不同的任务。

## 1.Master Thread

Master Thread是一个非常核心的后台线程，主要负责将缓冲池中的数据异步刷新到磁盘，保证数据的一致性，包括脏页的刷新、合并插入缓冲（INSERT BUFFER）、UNDO页的回收等。2.5节会详细地介绍各个版本下Master Thread的工作方式。

## 2.IO Thread

在InnoDB存储引擎中大量使用了AIO（Async IO）来处理写IO请求，这样可以极大提高数据库的性能。而IO Thread的工作主要是负责这些IO请求的回调（call back）处理。InnoDB 1.0版本之前共有4个IO Thread，分别是write、read、insert buffer和log IO thread。在Linux平台下，IO Thread的数量不能进行调整，但是在Windows平台下可以通过参数innodb_file_io_threads来增大IO Thread。从InnoDB 1.0.x版本开始，read thread和write thread分别增大到了4个，并且不再使用innodb_file_io_threads参数，而是分别使用innodb_read_io_threads和innodb_write_io_threads参数进行设置，如：

---

```
mysql＞SHOW VARIABLES LIKE'innodb_version'\G;

***************************1.row***************************

Variable_name:innodb_version

Value:1.0.6

1 row in set(0.00 sec)

mysql＞SHOW VARIABLES LIKE'innodb_%io_threads'\G;

***************************1.row***************************

Variable_name:innodb_read_io_threads

Value:4

***************************2.row***************************
```

Variable_name:innodb_write_io_threads

Value:4

2 rows in set(0.00 sec)

---

# 可以通过命令SHOW ENGINE INNODB STATUS来观察InnoDB中的IO Thread：

mysql：SHOW ENGINE INNODB STATUS\G;

***************************1.row***************************

Type:InnoDB

Name:

Status:

=====================================

100719 21:55:26 INNODB MONITOR OUTPUT

=====================================

Per second averages calculated from the last 36 seconds

……

--------

FILE I/O

--------

I/O thread 0 state:waiting for i/o request(insert buffer thread)

I/O thread 1 state:waiting for i/o request(log thread)

I/O thread 2 state:waiting for i/o request(read thread)

I/O thread 3 state:waiting for i/o request(read thread)

I/O thread 4 state:waiting for i/o request(read thread)

I/O thread 5 state:waiting for i/o request(read thread)

I/O thread 6 state:waiting for i/o request(write thread)

I/O thread 7 state:waiting for i/o request(write thread)

I/O thread 8 state:waiting for i/o request(write thread)

I/O thread 9 state:waiting for i/o request(write thread)

......

----------------------------

END OF INNODB MONITOR OUTPUT

============================

1 row in set(0.01 sec)

---

这里的IO Thread 0为insert buffer thread，IO Thread 1为log thread。之后就是根据参数innodb_read_io_threads及innodb_write_io_threads来设置的读写线程，并且读线程的ID总是小于写线程。

## 3.Purge Thread

事务被提交后，其所使用的undolog可能不再需要，因此需要PurgeThread来回收已经使用并分配的undo页。在InnoDB 1.1版本之前，purge操作仅在InnoDB存储引擎Master Thread中完成。而从InnoDB 1.1版本开始，purge操作可以独立到单独的线程中进行，以此来减轻Master Thread的工作，从而提高CPU的使用率以及提升存储引擎的性能。用户可以在MySQL数据库的配置文件中添加如下命令来启用独立的Purge Thread：

```
[mysqld]

innodb_purge_threads=1
```

在InnoDB 1.1版本中，即使将innodb_purge_threads设为大于1，InnoDB存储引擎启动时也会将其设为1，并在错误文件中出现如下类似的提示：

```
120529 22:54:16[Warning]option'innodb-purge-threads':unsigned value 4 adjusted to 1
```

从InnoDB 1.2版本开始，InnoDB支持多个Purge Thread，这样做的目的是为了进一步加快undo页的回收。同时由于Purge Thread需要离散地读取undo页，这样也能更进一步利用磁盘的随机读取性能。如用户可以设置4个Purge Thread：

```
mysql：SELECT VERSION()\G;

*************************1.row*************************
```

```
VERSION():5.6.6

1 row in set(0.00 sec)

mysql SHOW VARIABLES LIKE'innodb_purge_threads'\G;

***************************1.row***************************

Variable_name:innodb_purge_threads

Value:4

1 row in set(0.00 sec)
```

## 4.Page Cleaner Thread

Page Cleaner Thread是在InnoDB 1.2.x版本中引入的。其作用是将之前版本中脏页的刷新操作都放入到单独的线程中来完成。而其目的是为了减轻原Master Thread的工作及对于用户查询线程的阻塞，进一步提高InnoDB存储引擎的性能。

# 2.3.2 内存

## 1.缓冲池

InnoDB存储引擎是基于磁盘存储的，并将其中的记录按照页的方式进行管理。因此可将其视为基于磁盘的数据库系统（Disk-base Database）。在数据库系统中，由于CPU速度与磁盘速度之间的鸿沟，基于磁盘的数据库系统通常使用缓冲池技术来提高数据库的整体性能。

缓冲池简单来说就是一块内存区域，通过内存的速度来弥补磁盘速度较慢对数据库性能的影响。在数据库中进行读取页的操作，首先将从磁盘读到的页存放在缓冲池中，这个过程称为将页"FIX"在缓冲池中。下一次再读取相同的页时，首先判断该页是否在缓冲池中。若在缓冲池中，称该页在缓冲池中被命中，直接读取该页。否则，读取磁盘上的页。

对于数据库中页的修改操作，则首先修改在缓冲池中的页，然后再以一定的频率刷新到磁盘上。这里需要注意的是，页从缓冲池刷新回磁盘的操作并不是在每次页发生更新时触发，而是通过一种称为Checkpoint的机制刷新回磁盘。同样，这也是为了提高数据库的整体性能。

对于InnoDB存储引擎而言，其缓冲池的配置通过参数32位操作系统能管理的内存是有限的，最多只有3G，有一些开启了操作系统的PAE功能以此让32位操作系统能管理64GB的内存。但是现在数据库应用多运行于64位操作系统上，内存早已突破了8GB的限制。我现在使用的PC内存都已经达到512GB了，因此本书之后所讨论的内存管理的架构都将基于64位的操作系统。

由于InnoDB存储引擎的内存是通过参数innodb_buffer_pool_size来设置的。下面显示一台MySQL数据库服务器，其InnoDB存储引擎缓冲池的设置为15GB。

---

```
mysql>SHOW VARIABLES LIKE'innodb_buffer_pool_size'\G;

***************************1.row***************************

Variable_name:innodb_buffer_pool_size

Value:16106127360

1 row in set(0.00 sec)
```

---

具体来看，缓冲池中缓存的数据页类型有：索引页、数据页、undo页、插入缓冲（insert buffer）、自适应哈希索引（adaptive hash index）、InnoDB存储的锁信息（lock info）、数据字典信

数据字典（data dictionary）、插入缓冲、自适应哈希索引、锁信息（lock info）等，如图2-2所示，具体介绍了InnoDB存储引擎内存的结构情况。



图 2-2  InnoDB内存数据对象

从InnoDB 1.0.x版本开始，缓冲池的配置可以有多个缓冲池实例。每个页根据哈希值平均分配到不同缓冲池实例中。这样做的好处是减少数据库内部的资源竞争，增加数据库的并发处理能力。可以通过参数innodb_buffer_pool_instances来进行配置，该值默认为1。

```
mysql>SHOW VARIABLES LIKE'innodb_buffer_pool_instances'\G;
```

```
***************************1.row***************************
```

Variable_name:innodb_buffer_pool_instances

Value:1

1 row in set(0.00 sec)

---

从上面可以看出innodb_buffer_pool_instances的值为1，说明只有一个缓冲池实例。运行SHOW ENGINE INNODB STATUS命令查看缓冲池的使用情况：

---

mysql，SHOW ENGINE INNODB STATUS\G;

*****************************1.row***************************

Type:InnoDB

……

----------------------

INDIVIDUAL BUFFER POOL INFO

----------------------

---BUFFER POOL 0

Buffer pool size 65535

Free buffers 65451

Database pages 84

Old database pages 0

Modified db pages 0

Pending reads 0

Pending writes:LRU 0,flush list 0 single page 0

Pages made young 0,not young 0

0.00 youngs/s,0.00 non-youngs/s

Pages read 84,created 0,written 1

9.33 reads/s,0.00 creates/s,0.11 writes/s

Buffer pool hit rate 764/1000,young-making rate 0/1000 not 0/1000

Pages read ahead 0.00/s,evicted without access 0.00/s,Random read ahead 0.00/s

LRU len:84,unzip_LRU len:0

I/O sum[0]:cur[0],unzip sum[0]:cur[0]

```
---BUFFER POOL 1

Buffer pool size 65536

Free buffers 65473

Database pages 63

Old database pages 0

Modified db pages 0

Pending reads 0

Pending writes:LRU 0,flush list 0 single page 0

Pages made young 0,not young 0

0.00 youngs/s,0.00 non-youngs/s

Pages read 63,created 0,written 0

7.00 reads/s,0.00 creates/s,0.00 writes/s

Buffer pool hit rate 500/1000,young-making rate 0/1000 not 0/1000

Pages read ahead 0.00/s,evicted without access 0.00/s,Random read ahead 0.00/s

LRU len:63,unzip_LRU len:0

I/O sum[0]:cur[0],unzip sum[0]:cur[0]
```

由于参数innodb_buffer_pool_instances默认值为2，因此在上述状态文件中可以看到两部分的BUFFER POOL信息，为了节省空间，笔者省略了另外一部分（---BUFFER POOL 0）的显示，其信息和上面的大致相同。

从MySQL 5.6版本开始，还可以通过information_schema架构下的表INNODB_BUFFER_POOL_STATS来观察缓冲池的运行状态，如下所示的SQL语句所示：

```
mysql>SELECT POOL_ID,POOL_SIZE,

-  FREE_BUFFERS,DATABASE_PAGES

-  FROM INNODB_BUFFER_POOL_STATS\G;

***************************1.row***************************

POOL_ID:0

POOL_SIZE:65535
```

```
FREE_BUFFERS:65451

DATABASE_PAGES:84

***************************2.row***************************

POOL_ID:1

POOL_SIZE:65536

FREE_BUFFERS:65473

DATABASE_PAGES:63
```

---

## 2.LRU List、Free List、Flush List

通常来说，数据库中的缓冲池是通过管理的。在了解管理之前，读者必须知道InnoDB存储引擎是怎么对缓冲池进行管理的，即怎么对这么大的内存区域进行管理。

通常来说，数据库中的缓冲池是通过LRU（Latest Recent Used，最近最少使用）算法来进行管理的。即最频繁使用的页在LRU列表的前端，而最少使用的页在LRU列表的尾端。当缓冲池不能存放新读取到的页时，将首先释放LRU列表中尾端的页。

在InnoDB存储引擎中，缓冲池中页的大小默认为16KB，同样使用LRU算法对缓冲池进行管理。稍有不同的是InnoDB存储引擎对传统的LRU算法做了一些优化。在InnoDB的存储引擎中，LRU列表中还加入了midpoint位置。新读取到的页，虽然是最新访问的页，但并不是直接放入到LRU列表的首部，而是放入到LRU列表的midpoint位置。这个算法在InnoDB存储引擎下称为midpoint insertion strategy。在默认配置下，该位置在LRU列表长度的5/8处。midpoint位置可由参数innodb_old_blocks_pct控制，如：

---

```
mysql>SHOW VARIABLES LIKE'innodb_old_blocks_pct'\G;

***************************1.row***************************

Variable_name:innodb_old_blocks_pct

Value:37

1 row in set(0.00 sec)
```

---

从上面的例子可以看到，参数innodb_old_blocks_pct默认值为37，表示新读取的页插入到LRU列表尾端的37%的位置（差不多3/8的位置）。在InnoDB存储引擎中，把midpoint之后的列表称为old列表，之前的列表称为new列表。可以简单地理解为new列表中的页都是最为活跃的热点数据。

□□□□□□□□□□□LRU□□□□□□□□□□□□□□□□LRU□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□LRU□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□LRU□□□□□□□□□□□□□□□□□□□□□□□□□□□□LRU□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□

□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□LRU□□□□□□□□□□□innodb_old_blocks_time□□□□□□□□□□mid□□□□□□□□□□□□□□□□□□LRU□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□LRU□□□□□□□□□□□□□□

---

```
mysql□SET GLOBAL innodb_old_blocks_time=1000;

Query OK,0 rows affected(0.00 sec)

#data or index scan operation

……

mysql□SET GLOBAL innodb_old_blocks_time=0;

Query OK,0 rows affected(0.00 sec)
```

---

□□□□□□□□□□□□□□□□□63%□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

```
mysql□SET GLOBAL innodb_old_blocks_pct=20;

Query OK,0 rows affected(0.00 sec)
```

---

LRU□□□□□□□□□□□□□□□□□□□□□□□□□□□□LRU□□□□□□□□□□□□□□□□□□□□□□□□□Free□□□□□□□□□□□□□□□□□□□□□□□□Free□□□□□□□□□□□□□□□□□□□□□□□□□Free□□□□□□□□□LRU□□□□□□□□□□LRU□□□□□LRU□□□□□□□□□□□□□□□□□□□□□□LRU□□□old□□□□□new□□□□□□□□□□□□□page made young□□□□innodb_old_blocks_time□□□□□□□□□□□□old□□□□□new□□□□□□□page not made young□□□□□□□SHOW ENGINE INNODB STATUS□□□LRU□□□Free□□□□□□□□□□□□□

---

```
mysql□SHOW ENGINE INNODB STATUS\G;
```

```
***************************1.row**************************

Type:InnoDB

Name:

Status:

================================

120725 22:04:25 INNODB MONITOR OUTPUT

================================

Per second averages calculated from the last 24 seconds

……

Buffer pool size 327679

Free buffers 0

Database pages 307717

Old database pages 113570

Modified db pages 24673

Pending reads 0

Pending writes:LRU 0,flush list 0,single page 0

Pages made young 6448526,not young 0

48.75 youngs/s,0.00 non-youngs/s

Pages read 5354420,created 239625,written 3486063

55.68 reads/s,81.74 creates/s,955.88 writes/s

Buffer pool hit rate 1000/1000,young-making rate 0/1000 not 0/1000

……
```

通过命令SHOW ENGINE INNODB STATUS可以观察到Buffer pool size共有327 679，也就是327679*16K，共计5GB的缓冲池。Free buffers表示当前Free列表中页的数量，Database pages表示LRU列表中页的数量。可能的情况是，Free buffers与Database pages的数量之和不等于Buffer pool size。正如图2-2中所示的那样，因为缓冲池中的页还可能会被分配给自适应哈希索引、Lock信息、Insert Buffer等页，而这部分页不需要LRU算法进行维护，因此不存在于LRU列表中。

pages made young（因为LRU列表中读取到热端的页面数）、页面由于未单独访问（由于innodb_old_blocks_time的设置而导致的not young）、0）youngs/s、non-youngs/s等信息，这个信息同样非常重要。这里有一个概念——Buffer pool hit rate，表示缓冲池的命中率，这个例子中为100%，说明缓冲池运行状态非常良好。通常该值不应该小于95%。若发生Buffer pool hit rate的值小于95%这种情况，用户需要观察是否是由于全表扫描引起的LRU列表被污染的问题。

注意 　执行命令SHOW ENGINE INNODB STATUS显示的不是当前的状态，而是过去某个时间范围内InnoDB存储引擎的状态。因为在上述例子中，Per second averages calculated from the last 24 seconds代表的信息为过去24秒内的数据库状态。

在InnoDB 1.2版本中，还加入了表INNODB_BUFFER_POOL_STATS来观察缓冲池的运行状态，如：

---

```
mysql＞SELECT POOL_ID,HIT_RATE,

-＞PAGES_MADE_YOUNG,PAGES_NOT_MADE_YOUNG

-＞FROM information_schema.INNODB_BUFFER_POOL_STATS\G;

***************************1.row***************************

POOL_ID:0

HIT_RATE:980

PAGES_MADE_YOUNG:450

PAGES_NOT_MADE_YOUNG:0
```

---

同时用户还可以通过表INNODB_BUFFER_PAGE_LRU来观察每个LRU列表中每个页的具体信息，例如通过下面的语句可以看到缓冲池LRU列表中SPACE为1的表的页类型：

---

```
mysql＞SELECT TABLE_NAME,SPACE,PAGE_NUMBER,PAGE_TYPE

-＞FROM INNODB_BUFFER_PAGE_LRU WHERE SPACE=1;

+------------+-------+------------+------------------+

|TABLE_NAME|SPACE|PAGE_NUMBER|PAGE_TYPE|

+------------+-------+------------+------------------+

|NULL|1|0|FILE_SPACE_HEADER|
```

```
|NULL|1|1|IBUF_BITMAP|

|NULL|1|2|INODE|

|test/t|1|3|INDEX|

+------------+-------+------------+------------------+

4 rows in set(0.00 sec)
```

---

InnoDB从存储引擎1.0.x版本开始支持压缩页的功能，即将原本16KB的页压缩为1KB、2KB、4KB和8KB。而由于页的大小发生了变化，LRU列表也有了些许的改变。对于非16KB的页，是通过unzip_LRU列表进行管理的。通过命令SHOW ENGINE INNODB STATUS可以观察到如下内容：

---

```
mysql：SHOW ENGINE INNODB STATUS\G;

……

Buffer pool hit rate 999/1000,young-making rate 0/1000 not 0/1000

Pages read ahead 0.00/s,evicted without access 0.00/s,Random read ahead 0.00/s

LRU len:1539,unzip_LRU len:156

I/O sum[0]:cur[0],unzip sum[0]:cur[0]

……
```

---

可以看到LRU列表中一共有1539个页，而unzip_LRU列表中有156个页。这里需要注意的是，LRU中的页包含了unzip_LRU列表中的页。

对于压缩页的表，每个表的压缩比率可能各不相同。可能有的表页大小为8KB，有的表页大小为2KB。所以unzip_LRU是区别对待不同压缩页大小的页。其次，要将压缩页分配到unzip_LRU列表中时。首先会检查是否有可用的空闲空间，若没有，则从LRU列表中淘汰页。以从缓冲池中申请分配一个4KB大小的内存区域为例。

1、检查4KB的unzip_LRU列表，检查是否有可用的空闲页；

2、若有，则直接使用；

3、否则，检查8KB的unzip_LRU列表；

4、如果没有符合条件的内存块（2和4KB的空闲内存块4KB的unzip_LRU链表。

5、遍历当前大小为分的LRU链表，找到一个16KB的内存块，分成1个8KB和两个2、4KB的内存块。后两者放入unzip_LRU链表中。

可以通过查询information_schema数据库下的INNODB_BUFFER_PAGE_LRU来观察unzip_LRU链表中的页：

mysql>SELECT

-> TABLE_NAME,SPACE,PAGE_NUMBER,COMPRESSED_SIZE

-> FROM INNODB_BUFFER_PAGE_LRU

-> WHERE COMPRESSED_SIZE<>0;

+------------+-------+-------------+-----------------+

|TABLE_NAME|SPACE|PAGE_NUMBER|COMPRESSED_SIZE|

+------------+-------+-------------+-----------------+

|sbtest/t|9|134|8192|

|sbtest/t|9|135|8192|

|sbtest/t|9|96|8192|

|sbtest/t|9|136|8192|

|sbtest/t|9|32|8192|

|sbtest/t|9|97|8192|

|sbtest/t|9|137|8192|

|sbtest/t|9|98|8192|

……

在LRU列表中的页被修改后，称该页为脏页（dirty page），即缓冲池中的页和磁盘上的页的数据产生了不一致。这时数据库会通过CHECKPOINT机制将脏页刷新回磁盘，而Flush列表中的页即为脏页列表。需要注意的是，脏页既存在于LRU列表中，也存在于Flush列表中。LRU列表用来管理缓冲池中页的可用性，Flush列表用来管理将页刷新回磁盘，二者互不影响。

同LRU列表一样，Flush列表中的页也可以通过SHOW ENGINE INNODB STATUS来查看，前面例子中的Modified db pages 24673就显示了脏页的数量。

information_schema架构下还有一张INNODB_BUFFER_PAGE_LRU表，该表能够显示缓冲池中每个页的具体信息，如在LRU列表中的位置等。下面通过查询INNODB_BUFFER_PAGE_LRU表显示了各个页的类型，其中OLDEST_MODIFICATION大于0的SQL语句如下所示。

---

mysql＞SELECT TABLE_NAME,SPACE,PAGE_NUMBER,PAGE_TYPE

-＞FROM INNODB_BUFFER_PAGE_LRU

-＞WHERE OLDEST_MODIFICATION＞0;

+------------+-------+-------------+-------------------+

|TABLE_NAME|SPACE|PAGE_NUMBER|PAGE_TYPE|

+------------+-------+-------------+-------------------+

|NULL|0|56|SYSTEM|

|NULL|0|0|FILE_SPACE_HEADER|

|test/t|1|3|INDEX|

|NULL|0|320|INODE|

|NULL|0|325|UNDO_LOG|

+------------+-------+-------------+-------------------+

5 rows in set(0.00 sec)

---

从上表中可以看到，有5个页的类型为脏页，其中页的TABLE_NAME为NULL的表示该页属于系统表空间。

## 3.重做日志缓冲

如图2-2所示，InnoDB存储引擎的内存区域除了有缓冲池外，还有重做日志缓冲（redo log buffer）。InnoDB存储引擎首先将重做日志信息先放入这个缓冲区，然后按一定频率将其刷新到重做日志文件。重做日志缓冲一般不需要设置得很大，因为一般情况下每一秒钟会将重做日志缓冲刷新到日志文件，因此用户只需要保证每秒产生的事务量在这个缓冲大小之内即可。该值可由配置参数innodb_log_buffer_size控制，默认为8MB。

---

mysql＞SHOW VARIABLES LIKE'innodb_log_buffer_size'\G;

***************************1.row***************************

```
Variable_name:innodb_log_buffer_size

Value:8388608

1 row in set(0.00 sec)
```

默认大小即为8MB，一般来说该值并不需要调整得太大。因为每一秒钟都会将重做日志缓冲刷新到日志文件，因此用户只需要保证每秒产生的事务量在这个缓冲大小之内即可。该值的设置取决于对重做日志缓冲刷新操作，有以下三种情况会将重做日志缓冲中的内容刷新到外部磁盘的重做日志文件中。

☐Master Thread每一秒将重做日志缓冲刷新到重做日志文件。

☐每个事务提交时会将重做日志缓冲刷新到重做日志文件。

☐当重做日志缓冲池剩余空间小于1/2时，重做日志缓冲刷新到重做日志文件。

## 4.额外的内存池

在对数据库的学习过程中，DBA可能会忽视了这部分内存的使用。额外的内存池通常被DBA忽略，他们认为该值并不十分重要，事实恰恰相反，该值同样十分重要。在InnoDB存储引擎中，对内存的管理是通过一种称为内存堆（heap）的方式进行的。在对一些数据结构本身的内存进行分配时，需要从额外的内存池中进行申请，当该区域的内存不够时，会从缓冲池中进行申请。例如，分配了innodb_buffer_pool，但是每个缓冲池中的帧缓冲（frame buffer）还有对应的缓冲控制对象（buffer control block），这些对象记录了一些诸如LRU、锁、等待等信息，而这个对象的内存需要从额外内存池中申请。因此，在申请了很大的InnoDB缓冲池时，也应考虑相应地增加这个值。

# 2.4 Checkpoint技术

在数据库发展的早期阶段，还是用CPU的处理能力非常有限的情况下，如果一个数据库用户执行了一条DML语句（如Update、Delete），那么数据库系统需要将这些修改立即同步到磁盘上的数据文件中，才能算执行成功。

但是这样处理会有两个很致命的问题，首先磁盘写入的速度是非常慢的，如果每一次修改都要同步写入磁盘，那么数据库的性能将会非常低下。其次，如果在修改的过程中出现了宕机等情况，那么数据库系统就需要一种机制来保证数据的一致性。因此，数据库系统引入了日志（Write Ahead Log）机制，将这些修改先记录到日志文件中，然后再异步地将这些修改同步到磁盘上的数据文件中，这样就可以保证数据的持久性（ACID中D，Durability）和一致性了。

但是日志文件也不能无限增长，否则会占用大量的磁盘空间，因此数据库系统引入了检查点机制，定期地将内存中的修改同步到磁盘上的数据文件中，并记录一个检查点，这样就可以清理掉检查点之前的日志文件了。

□ 检查点机制可以有效地减少日志文件的大小

□ 检查点机制可以加快数据库

数据库系统在进行检查点操作时，需要将内存中的所有修改同步到磁盘上的数据文件中，这个过程可能会非常耗时，尤其是在数据量非常大的情况下。例如一个大型数据库的内存缓冲区达到3TB，MySQL数据库的缓冲池达到3 TB，一个大型的Oracle Exadata数据库的内存缓冲区达到2 TB，如果要将这些修改全部同步到磁盘上，那么需要花费大量的时间。

因此，数据库系统在进行检查点操作时，需要尽可能地减少对数据库性能的影响，这就需要DBA、SA等专业人员对数据库系统进行合理的配置和优化，才能保证数据库系统的高性能和高可用性。

检查点机制是数据库系统中一个非常重要的机制，它可以有效地保证数据的持久性和一致性，同时还可以减少日志文件的大小，提高数据库系统的性能。

那么Checkpoint机制具体是如何实现的呢？它又有哪些类型呢？

□ 检查点机制的实现原理

❑□□□□□□□□□□□□□□□□□

❑□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□Checkpoint□□□□□□□□□□□□□□□□□□□□□□□Checkpoint□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□LRU□□□□□□□□□□□□□□□□□□□□□□□□□□□Checkpoint□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Checkpoint□□□□□□□□□□□□□□□□□□□□□□□□□

□□InnoDB□□□□□□□□□□□□LSN□Log Sequence Number□□□□□□□□□□□LSN□8□□□□□□□□□□□□□□□□□□LSN□□□□□□□□□□LSN□Checkpoint□□LSN□□□□□□□□SHOW ENGINE INNODB STATUS□□□□□

---

```
mysql□SHOW ENGINE INNODB STATUS\G;

......

---

LOG

---

Log sequence number 92561351052

Log flushed up to   92561351052

Last checkpoint at  92561351052

......
```

---

□InnoDB□□□□□□□Checkpoint□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Checkpoint□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Checkpoint□□InnoDB□□□□□□□□□□□□Checkpoint□□□□□

☐Sharp Checkpoint

☐Fuzzy Checkpoint

Sharp Checkpoint发生在数据库关闭时将所有的脏页都刷新回磁盘，这是默认的工作方式，即参数innodb_fast_shutdown=1。

但是若在数据库运行时也使用Sharp Checkpoint，那么数据库的可用性就会受到很大的影响。故在InnoDB存储引擎内部使用Fuzzy Checkpoint进行页的刷新，即只刷新一部分脏页，而不是刷新所有的脏页回磁盘。

这里笔者进行了概括，在InnoDB存储引擎中可能发生如下几种情况的Fuzzy Checkpoint：

☐Master Thread Checkpoint

☐FLUSH_LRU_LIST Checkpoint

☐Async/Sync Flush Checkpoint

☐Dirty Page too much Checkpoint

对于Master Thread中发生的Checkpoint，差不多以每秒或每十秒的速度从缓冲池的脏页列表中刷新一定比例的页回磁盘。这个过程是异步的，即此时InnoDB存储引擎可以进行其他的操作，用户查询线程不会阻塞。

FLUSH_LRU_LIST Checkpoint是因为InnoDB存储引擎需要保证LRU列表中需要有差不多100个空闲页可供使用。在InnoDB1.1.x版本之前，需要检查LRU列表中是否有足够的可用空间操作发生在用户查询线程中，显然这会阻塞用户的查询操作。倘若没有100个可用空闲页，那么InnoDB存储引擎会将LRU列表尾端的页移除。如果这些页中有脏页，那么需要进行Checkpoint，而这些页是来自LRU列表的，因此称为FLUSH_LRU_LIST Checkpoint。

而从MySQL 5.6版本，也就是InnoDB1.2.x版本开始，这个检查被放在了一个单独的Page Cleaner线程中进行，并且用户可以通过参数innodb_lru_scan_depth控制LRU列表中可用页的数量，该值默认为1024，如：

mysql>SHOW VARIABLES LIKE'innodb_lru_scan_depth'\G;

***************************1.row***************************

Variable_name:innodb_lru_scan_depth

Value:1024

1 row in set(0.00 sec)

---

Async/Sync Flush Checkpoint指的是重做日志不可用的情况，这时需要强制刷新页回磁盘，而需要刷新的页是脏页列表中选取的。因此这种情况一般发生在重做日志很大的情况下。在这种情况下，重做日志可重用的数量需要判断的两个值：当前写入重做日志的LSN（即redo_lsn）和最后一次检查点的LSN（即checkpoint_lsn）它们的差值。

---

checkpoint_age=redo_lsn-checkpoint_lsn

---

## 又定义了两个变量：

---

async_water_mark=75%*total_redo_log_file_size

sync_water_mark=90%*total_redo_log_file_size

---

若每个重做日志文件的大小为1GB，并且定义了两个重做日志文件，则重做日志文件的总大小为2GB。那么async_water_mark=1.5GB，sync_water_mark=1.8GB。则：

□当checkpoint_age＜async_water_mark时，不需要刷新任何脏页到磁盘；

□当async_water_mark＜checkpoint_age＜sync_water_mark时触发Async Flush，从Flush列表中刷新足够的脏页回磁盘，使得刷新后满足checkpoint_age＜async_water_mark；

□checkpoint_age＞sync_water_mark这种情况一般很少发生，除非设置的重做日志文件太小，并且在进行类似LOAD DATA的BULK INSERT操作。此时触发Sync Flush操作，从Flush列表中刷新足够的脏页回磁盘，使得刷新后满足checkpoint_age＜async_water_mark。

可见，Async/Sync Flush Checkpoint是为了保证重做日志的循环使用的可用性。在InnoDB 1.2.x版本之前，Async Flush Checkpoint会阻塞发现问题的用户查询线程，而Sync Flush Checkpoint会阻塞所有的用户查询线程，并且等待脏页刷新完成。从

InnoDB 1.2.x版本之前——也就是MySQL 5.6版本之前，对于刷新脏页的操作是由Master Thread（或Page Cleaner Thread）来负责的。

MySQL中几乎所有的物理刷新操作，包括Flush操作（比如对LRU中脏页进行Checkpoint）或者是脏页过多时的Async/Sync Flush机制，在InnoSQL里面都可以通过命令SHOW ENGINE INNODB STATUS来进行查看。

---

mysql〉SHOW ENGINE INNODB STATUS\G;

***************************1.row***************************

Type:InnoDB

……

LRU len:112902,unzip_LRU len:0

I/O sum[0]:cur[0],unzip sum[0]:cur[0]

Async Flush:0,Sync Flush:0,LRU List Flush:0,Flush List Flush:111736

……

1 row in set(0.01 sec)

---

那么具体的刷新机制是什么？InnoDB刷新脏页的机制，其实就是围绕以下9点来进行的

第一点：Checkpoint机制（Dirty Page too much）。通过这个机制，可以控制InnoDB中的脏页。一旦Checkpoint不及时，脏页就会过多。这个机制可以通过innodb_max_dirty_pages_pct来控制。

---

mysql〉SHOW VARIABLES LIKE'innodb_max_dirty_pages_pct'\G;

***************************1.row***************************

Variable_name:innodb_max_dirty_pages_pct

Value:75

1 row in set(0.00 sec)

---

innodb_max_dirty_pages_pct值为75表示，当缓冲池中脏页的数量占据75%时，强制进行Checkpoint，刷新一部分的脏页到磁盘。在InnoDB 1.0.x版本之前，这个参数的默认

90□□□□□□□□75□

## 2.5 Master Thread工作方式

在2.3节中我们知道，InnoDB存储引擎的主要工作都是在一个单独的后台线程Master Thread中完成的，这一节将具体解释该线程的具体实现以及该线程可能存在的问题。

## 2.5.1 InnoDB 1.0.x版本之前的Master Thread

Master Thread具有最高的线程优先级别。其内部由多个循环（loop）组成：主循环（loop）、后台循环（backgroup loop）、刷新循环（flush loop）、暂停循环（suspend loop）。Master Thread会根据数据库运行的状态在loop、background loop、flush loop和suspendloop中进行切换。

Loop被称为主循环，因为大多数的操作是在这个循环中，其中有两大部分的操作——每秒钟的操作和每10秒钟的操作。伪代码如下：

```
void master_thread(){

loop：

for(int i=0;i＜10;i++){

do thing once per second

sleep 1 second if necessary

}

do things once per ten seconds

goto loop;

}
```

可以看到，loop循环通过thread sleep来实现，这意味着所谓的每秒一次或每10秒一次的操作是不精确的。在负载很大的情况下可能会有延迟（delay），只能说大概在这个频率下。当然，InnoDB源代码中还通过了其他的方法来尽量保证这个频率。

每秒一次的操作包括：

□日志缓冲刷新到磁盘，即使这个事务还没有提交（总

□□□□□□□□□□□□□

□□□□□□100□InnoDB□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□background loop□□□□□□

□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□commit□□□□□□□□□□□

□□□□□□□□Insert Buffer□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□IO□□□□□□5□□□□□□5□□InnoDB□□□□□IO□□□□□□□□□□□□□□□□□□□□

□□□□□□100□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□buf_get_modified_ratio_pct□□□□□□□□□□□□□innodb_max_dirty_pages_pct□□□□□□□□□□90□□□90%□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□100□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

```
void master_thread(){

goto loop;

loop□

for(int i=0;i□10;i++){

thread_sleep(1)//sleep 1 second

do log buffer flush to disk

if(last_one_second_ios□5)

do merge at most 5 insert buffer

if(buf_get_modified_ratio_pct□innodb_max_dirty_pages_pct)

do buffer pool flush 100 dirty page

if(no user activity)

goto backgroud loop

}

do things once per ten seconds
```

```
background loop:

do something

goto loop:

}
```

□每隔10秒会执行的操作有如下几个

□刷新100个脏页到磁盘（可能的情况下）

□合并至多5个插入缓冲（总是）

□刷新日志缓冲到磁盘（总是）

□删除无用的Undo页（总是）

□产生100次或者10次以上的主线程睡眠

在磁盘刷新时，InnoDB存储引擎会先判断过去10秒之内磁盘的IO操作是否小于200次，如果是，InnoDB存储引擎认为当前有足够的磁盘IO操作能力，所以将100个脏页刷新到磁盘。然后，InnoDB存储引擎会判断缓冲池中脏页的比例是否超过了配置文件中的参数，如果超过，InnoDB存储引擎同样会将脏页刷新回磁盘，刷新数量同上述策略一致。

接着，InnoDB存储引擎会进行一步执行full purge操作，即删除无用的Undo页。对表进行update、delete这类操作时，原先的行被标记为删除，但是因为一致性读（consistent read）的关系，需要保留这些行版本的信息。但在full purge过程中，InnoDB存储引擎会判断当前事务系统中已被删除的行是否可以删除，比如有时候可能还有查询操作需要读取之前版本的undo信息，如果可以删除，InnoDB存储引擎会立即将其删除。从源代码中可以发现，InnoDB存储引擎在执行full purge操作时，每次最多尝试回收20个undo页。

然后，InnoDB存储引擎会判断缓冲池中脏页的比例（buf_get_modified_ratio_pct），如果有超过70%的脏页，则刷新100个脏页到磁盘，如果脏页的比例小于70%，则只需刷新10%的脏页到磁盘。

现在我们回到主循环的伪代码，讨论main loop中每10秒的操作，如下所示。

```
void master_thread(){

goto loop;

loop:

for(int i=0;i<10;i++){

thread_sleep(1)//sleep 1 second

do log buffer flush to disk

if(last_one_second_ios<5)

do merge at most 5 insert buffer

if(buf_get_modified_ratio_pct>innodb_max_dirty_pages_pct)

do buffer pool flush 100 dirty page

if(no user activity)

goto backgroud loop

}

if(last_ten_second_ios<200)

do buffer pool flush 100 dirty page

do merge at most 5 insert buffer

do log buffer flush to disk

do full purge

if(buf_get_modified_ratio_pct>70%)

do buffer pool flush 100 dirty page

else

buffer pool flush 10 dirty page

goto loop

background loop:

do something

goto loop:

}
```

最后的background loop会在数据库关闭时进行切换，这里可以简单地理解为在shutdown时，会进行的操作，在这里background loop主要完成的工作

□ 回收已经使用并分配的Undo页；□

□ 刷新20个脏页到磁盘；□

□ 产生一个检查点；□

□ 刷新最多100个脏页到磁盘。接着，master thread会跳转到flush loop中去。

在flush loop中完成的大部分工作后，InnoDB存储引擎会跳转到suspend__loop，将Master Thread挂起，等待事件的发生。若当前没有用户enable了的InnoDB存储引擎，则Master Thread总是处于挂起的状态。Master Thread的伪代码如下：

若用户启用（enable）了InnoDB存储引擎，Master Thread的伪代码如下所示：

---

```
void master_thread(){

goto loop;

loop：

for(int i=0;i<10;i++){

thread_sleep(1)//sleep 1 second

do log buffer flush to disk

if(last_one_second_ios<5)

do merge at most 5 insert buffer

if(buf_get_modified_ratio_pct>innodb_max_dirty_pages_pct)

do buffer pool flush 100 dirty page

if(no user activity)

goto backgroud loop

}

if(last_ten_second_ios<200)

do buffer pool flush 100 dirty page

do merge at most 5 insert buffer

do log buffer flush to disk

do full purge
```

```
if(buf_get_modified_ratio_pct□70%)

do buffer pool flush 100 dirty page

else

buffer pool flush 10 dirty page

goto loop

background loop:

do full purge

do merge 20 insert buffer

if not idle:

goto loop:

else:

goto flush loop

flush loop:

do buffer pool flush 100 dirty page

if(buf_get_modified_ratio_pct□innodb_max_dirty_pages_pct)

goto flush loop

goto suspend loop

suspend loop:

suspend_thread()

waiting event

goto loop;

}
```

# 2.5.2　InnoDB1.2.x□□□□□Master Thread

□□□□□1.0.x□□□□□□Master Thread□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□IO□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□hard coding□□□□□□□□□□□□□□□□□□□□□□□□□□SSD□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□IO□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□100□□□□□□□□□□20□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□100□□□□□□□□□□□□□□20□□□□□□□□□□□□Master Thread□□□□□"□□□□"□□□□□□□□□□□□□□□□□□□□□1□□□□□□□□100□□□□□□□20□□□□□□□□□□□□□□□□hard coding□Master Thread□□□□□□□□□100□□□□□□□20□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□insert buffer□□□□

□□□□□□□□□Google□□□□□Mark Callaghan□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□patch□□InnoDB□□□□□□□□□□□□□□□Google□patch□□□□□□□□□□□□□□□□□□□□□InnoDB Plugin□□InnoDB1.0.x□□□□□□□□□□□□innodb_io_capacity□□□□□□□□□IO□□□□□□□□□□200□□□□□□□□□□□□□□□□□□□innodb_io_capacity□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□innodb_io_capacity□□□5%□

□□□□□□□□□□□□□□□□□□□□□□innodb_io_capacity□

□□□□□□□□□SSD□□□□□□□□□□□□□□□RAID□□□□□□□□□□□□□□□□IO□□□□□□□□□□□□innodb_io_capacity□□□□□□□□□□□□□□□□□IO□□□□□□□□□

□□□□□□□□□□□innodb_max_dirty_pages_pct□□□□□□□□□□□InnoDB 1.0.x□□□□□□□□□□□□□90□□□□□□□□□□□□□90%□□□□□□□"□□"□□□□□InnoDB□□□□□□□□□□□□□□□□flush loop□□□□□□□□□□□□□□□□□innodb_max_dirty_pages_pct□□□□□100□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□20□10□□□□□□□□□□□□□□□□□□□□□innodb_max_dirty_pages_pct□□20□10□□□□□□□□□□□□□□□□□□□□□□□□□□□Google□□□□□□□□□□□□□□□□20□□□□□□□□□ [1]□□□□InnoDB 1.0.x

参数。该innodb_max_dirty_pages_pct值默认为75，即Google设置的80。该值越大，每次刷新的脏页数量越多，可能对IO会有影响。

InnoDB 1.0.x版本引入了一个新的参数innodb_adaptive_flushing（自适应地刷新），该值影响每秒刷新脏页的数量。原来的刷新规则是：脏页在缓冲池所占的比例小于innodb_max_dirty_pages_pct时，不刷新脏页；大于innodb_max_dirty_pages_pct时，刷新100个脏页。随着innodb_adaptive_flushing参数的引入，InnoDB存储引擎会通过一个名为buf_flush_get_desired_flush_rate的函数来判断需要刷新脏页最合适的数量。粗略地翻译buf_flush_get_desired_flush_rate通过判断产生redo log的速度来决定最合适的刷新脏页数量。因此，当脏页的比例小于innodb_max_dirty_pages_pct时，也会刷新一定量的脏页。

还有一个操作是合并插入缓冲（full purge）。合并插入缓冲并非在每次Master Thread中，而是在InnoDB 1.0.x版本之前，InnoDB存储引擎会进行合并插入缓冲的操作。同样，每次进行full purge操作时，每次回收Undo页的数量为20，InnoDB 1.0.x版本开始引入了参数innodb_purge_batch_size，该参数可以控制每次full purge回收的Undo页的数量。该参数的默认值为20，并可以动态地对其进行修改，具体如下：

---

```
mysql>SHOW VARIABLES LIKE'innodb_purge_batch_size'\G;

***************************1.row***************************

Variable_name:innodb_purge_batch_size

Value:20

mysql>SET GLOBAL innodb_purge_batch_size=50;

Query OK,0 rows affected(0.00 sec)
```

---

从以上代码中可以看到，在InnoDB 1.0.x版本之前，Master Thread中每秒和每10秒的操作同样还存在着一些问题，那就是

---

```
void master_thread(){

goto loop;

loop：

for(int i=0;i<10;i++){

thread_sleep(1)//sleep 1 second

do log buffer flush to disk
```

```
if(last_one_second_ios<5%innodb_io_capacity)

do merge 5%innodb_io_capacity insert buffer

if(buf_get_modified_ratio_pct>innodb_max_dirty_pages_pct)

do buffer pool flush 100%innodb_io_capacity dirty page

else if enable adaptive flush

do buffer pool flush desired amount dirty page

if(no user activity)

goto backgroud loop

}

if(last_ten_second_ios<innodb_io_capacity)

do buffer pool flush 100%innodb_io_capacity dirty page

do merge 5%innodb_io_capacity insert buffer

do log buffer flush to disk

do full purge

if(buf_get_modified_ratio_pct>70%)

do buffer pool flush 100%innodb_io_capacity dirty page

else

dobuffer pool flush 10%innodb_io_capacity dirty page

goto loop

background loop:

do full purge

do merge 100%innodb_io_capacity insert buffer

if not idle:

goto loop:

else:

goto flush loop

flush loop:

do buffer pool flush 100%innodb_io_capacity dirty page

if(buf_get_modified_ratio_pct>innodb_max_dirty_pages_pct)

go to flush loop
```

```
goto suspend loop

suspend loop:

suspend_thread()

waiting event

goto loop;

}
```

---

可以看到,因为InnoDB 1.0.x版本之前的版本使用硬编码的方式,因此Master Thread的问题是不能更改的。而InnoDB存储引擎希望用户知道目前的Master Thread工作的状态。

从InnoDB 1.0.x开始,可以通过SHOW ENGINE INNODB STATUS命令来查看当前Master Thread的状态信息,如下所示。

---

```
mysql>SHOW ENGINE INNODB STATUS\G;

***************************1.row***************************

Type:InnoDB

Name:

Status:

=================================

090921 14:24:56 INNODB MONITOR OUTPUT

=================================

Per second averages calculated from the last 6 seconds

----------

BACKGROUND THREAD

----------

srv_master_thread loops:45 1_second,45 sleeps,4 10_second,6 background,6 flush

srv_master_thread log flush and writes:45 log writes only:69

……
```

---

最大休眠时间不能超过45。累积的休眠次数（sleep）超过一定次数（45）或者休眠时间超过一定值（10），就切换到定时循环（4秒，比1:10。background loop循环6次，flush loop循环达到6次，因为线程休眠有一定随机性和不确定性，因此这里的值只是理论值，不同版本的MySQL略有差异。以下是一个实际的例子：

---

mysql□show engine innodb status\G;

***************************1.row***************************

Type:InnoDB

Name:

Status:

==================================

091009 10:14:34 INNODB MONITOR OUTPUT

==================================

Per second averages calculated from the last 42 seconds

----------

BACKGROUND THREAD

----------

srv_master_thread loops:2188 1_second,1537 sleeps,218 10_second,2 background,2 flush

srv_master_thread log flush and writes:1777 log writes only:5816

……

---

可以看到，主循环进行了2188次，每秒挂起（sleep）的操作进行了1537次，说明当前InnoDB的压力并不大，如果主循环进行1次的时间用不到1_second，InnoDB内部就会进行sleep操作来控制频率。感兴趣的读者可以用自己的机器观察一下这个值的变化情况。

[1]_本节所举例子引用自 http://code.google.com/p/google-mysql-tools/wiki/InnodbIoOltpDisk。

## 2.5.3　InnoDB 1.2.x版本的Master Thread

在InnoDB 1.2.x版本之前的Master Thread已经有了很大的优化，而在InnoDB 1.2.x版本中对Master Thread的实现又进行了一些优化。在InnoDB 1.2.x版本中，Master Thread的伪代码如下：

---

```
if InnoDB is idle

srv_master_do_idle_tasks();

else

srv_master_do_active_tasks();
```

---

其中srv_master_do_idle_tasks()就是之前版本中每10秒的操作，srv_master_do_active_tasks()处理的是之前每秒中的操作。同时对于刷新脏页的操作，从Master Thread线程分离到一个单独的Page Cleaner Thread，从而减轻了Master Thread的工作，同时进一步提高了系统的并发性。

# 2.6　InnoDB□□□□

InnoDB□□□□□□□□□□□□□□

□□□□□□（Insert Buffer）

□□□□（Double Write）

□□□□□□□□（Adaptive Hash Index）

□□□IO（Async IO）

□□□□□□（Flush Neighbor Page）

□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□

## 2.6.1　□□□□

### 1.Insert Buffer

Insert Buffer□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□Insert Buffer□□□□□□□□□□Insert Buffer□□□□□□□□□□□□□□□□□□□□□□

□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Primary Key□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□

```
CREATE TABLE t(

a INT AUTO_INCREMENT,

b VARCHAR(30),

PRIMARY KEY(a)

);
```

列（a）的值为空值（a列允许NULL），或者没有指定（AUTO_INCREMENT时）一般情况下，此时用户可以不需要将主键设置为列a，即可达到同样的效果。但是当用户未指定主键时，系统会自动选择一个非空的唯一索引作为主键。

当然，还有一种情况是，用户使用诸如UUID的方式来生成主键的值。这时用户生成的主键值一般是无序的，并非是顺序递增NULL的（同时不允许为空），这时就会产生类似于下面代码。

下面来看一个例子。假如表有如下的索引结构，其中存在一个辅助索引（secondary index），有两个字段，分别为b的。且用户对于列b的值是顺序递增插入的，如下面的SQL语句所示。

CREATE TABLE t(

a INT AUTO_INCREMENT,

b VARCHAR(30),

PRIMARY KEY(a),

key(b)

);

在上面创建的表中存在两个索引，一个是聚集索引，一个是辅助索引。对于列a而言，由于其是自增长的，因此向表中插入行记录时，页中的行记录按照主键a进行顺序存放。而辅助索引上列b的值相对杂乱无序，这时B+树的存放还是按照键值进行存放的。

但是由于随机读取，这种插入操作效率会很低。然而，这并不是因为存在随机读取这一情况，只是由于辅助索引树的离散性造成了这个"问题"的存在。

InnoDB存储引擎开创性地设计了Insert Buffer，对于非聚集索引的插入或更新操作，不是每一次直接插入到索引页中，而是先判断插入的非聚集索引页是否在缓冲池中，若在，则直接插入；若不在，则先放入到一个Insert Buffer对象中，好似欺骗。数据库这个非聚集的索引已经插到叶子节点，而实际并没有，只是存放在另一个位置。然后再以一定的频率和情况进行Insert Buffer和辅助索引页子节点的merge（合并）操作，这时通常能将多个插入合并到一个操作中（因为在一个索引页中），这就大大提高了对于非聚集索引插入的性能。

然而Insert Buffer的使用需要同时满足以下两个条件：

□索引是辅助索引（secondary index）；

□索引不是唯一（unique）的。

这两个条件都满足的时候，InnoDB存储引擎会使用Insert Buffer，这样就能提高插入操作的性能了。不过这里还有一个问题，在写密集的情况下，插入缓冲会占用过多的Insert Buffer Buffer内存，默认MySQL最大可以占用的缓冲区内存为Insert Buffer的内存为缓冲池内存的一半，因此用户在一个写密集的应用下，可以考虑减小缓冲池的内存。

这两个条件都满足的时候，数据库的应用程序会对辅助索引的插入操作进行优化，这样就能提高插入操作的性能了，这就是Insert Buffer的设计目的。

用户可以通过命令SHOW ENGINE INNODB STATUS来查看插入缓冲的信息：

```
mysql SHOW ENGINE INNODB STATUS\G;

***************************1.row***************************

Type:InnoDB

Name:

Status:

=====================================

100727 22:21:48 INNODB MONITOR OUTPUT

=====================================

Per second averages calculated from the last 44 seconds

……

-------------------------------------

INSERT BUFFER AND ADAPTIVE HASH INDEX

-------------------------------------

Ibuf:size 7545,free list len 3790,seg size 11336,

8075308 inserts,7540969 merged recs,2246304 merges

……

----------------------------

END OF INNODB MONITOR OUTPUT

==========================
```

---

seg size显示了当前Insert Buffer的大小，为11336×16KB，大约为177MB；free list len代表了空闲列表的长度；size代表了已经合并记录页的数量。而黑体部分的第2行可能是用户最关心的，因为它显示了插入性能的提高。Inserts代表了插入的记录数；merged recs代表了合并的插入记录数量；merges代表合并的次数，也就是实际读取页的次数。merges:merged recs大约为1:3，代表了插入缓冲将对于非聚集索引页的离散IO逻辑请求大约降低了2/3。

正如前面所说的，目前Insert Buffer存在一个问题是：在写密集的情况下，插入缓冲会占用过多的缓冲池内存（innodb_buffer_pool），默认最大可以占用到1/2的缓冲池内存。以下代码显示了InnoDB源代码中对于insert buffer的初始化操作：

---

/**Buffer pool size per the maximum insert buffer size*/

#define IBUF_POOL_SIZE_PER_MAX_SIZE 2

ibuf-□max_size=buf_pool_get_curr_size()/UNIV_PAGE_SIZE

/IBUF_POOL_SIZE_PER_MAX_SIZE;

---

对于写密集的情况，插入缓冲会占用过多的缓冲池内存，因此在默认情况下并不适合。假设某个应用的操作全部是插入操作，若缓冲池的大小为100MB，则Insert Buffer最多可以使用50MB的缓冲池内存，这对于用户来说并不希望看到的。上述的代码用来判断Insert Buffer的大小，其中buf_pool_get_curr_size() 用来得到缓冲池的大小。因此当Percona发布一些patch来修正插入缓冲占用太多缓冲池内存的情况，具体可以到Percona官网进行查找。简单来说，修改IBUF_POOL_SIZE_PER_MAX_SIZE就可以对插入缓冲的大小进行控制。比如将IBUF_POOL_SIZE_PER_MAX_SIZE改为3，则最大只能使用1/3的缓冲池内存。

## 2.Change Buffer

InnoDB从1.0.x版本开始引入了Change Buffer，可将其视为Insert Buffer的升级。从这个版本开始，InnoDB存储引擎可以对DML操作——INSERT、DELETE、UPDATE都进行缓冲，他们分别是：Insert Buffer、Delete Buffer、Purge buffer。

当然和之前Insert Buffer一样，Change Buffer适用的对象依然是非唯一的辅助索引。

对一条记录进行UPDATE操作可能分为两个过程：

□将记录标记为已删除；

□真正将记录删除。

称为Delete Buffer，对UPDATE操作的第一个过程（标记为已删除）。Purge Buffer对UPDATE操作的第二个过程（真正将记录删除）。同时，InnoDB提供了参数innodb_change_buffering，用来开启各种Buffer的选项。该参数可选值为inserts、deletes、purges、changes、all、none。inserts、deletes、purges就是前面讨论过的三种情况。changes表示启用inserts和deletes，all表示启用所有，none表示都不启用。该参数默认值为all。

从InnoDB 1.2.x版本开始，可以通过参数innodb_change_buffer_max_size来控制Change Buffer最大使用内存的数量：

---

```
mysql>SHOW VARIABLES LIKE'innodb_change_buffer_max_size'\G;

***************************1.row***************************

Variable_name:innodb_change_buffer_max_size

Value:25

1 row in set(0.00 sec)
```

---

innodb_change_buffer_max_size值默认为25，表示最多使用1/4的缓冲池内存空间。而需要注意的是，该参数的最大有效值为50。

在MySQL 5.5版本中通过命令SHOW ENGINE INNODB STATUS，可以观察到类似如下的内容：

---

```
mysql>SHOW ENGINE INNODB STATUS\G;

***************************1.row***************************

Type:InnoDB

……

-------------------------------------

INSERT BUFFER AND ADAPTIVE HASH INDEX

-------------------------------------

Ibuf:size 1,free list len 34397,seg size 34399,10875 merges

merged operations:

insert 20462,delete mark 20158,delete 4215
```

```
discarded operations:

insert 0,delete mark 0,delete 0
```

……

---

可以看到，这里分别显示了merged operations和discarded operation，并记录了Change Buffer中每个操作的次数。insert表示Insert Buffer；delete mark表示Delete Buffer；delete表示Purge Buffer；discarded operations表示当Change Buffer发生merge时，表已经被删除，此时就无须再将记录merge到辅助索引中了。

## 3.Insert Buffer的内部实现

通过前面的介绍，读者应该已经对Insert Buffer的使用、作用有了一定的了解。也许读者会有这样的疑问，Insert Buffer的数据结构到底是怎样的？为什么叫Insert Buffer？

读者有这样的疑问是很正常的，因为我们之前并没有提到Insert Buffer的数据结构是一棵B+树。在MySQL 4.1之前的版本中每张表有一棵Insert Buffer B+树。而在现在的版本中，全局只有一棵Insert Buffer B+树，负责对所有的表的辅助索引进行Insert Buffer。而这棵B+树存放在共享表空间中，默认也就是ibdata1中。因此，试图通过独立表空间ibd文件恢复表中数据时，往往会导致CHECK TABLE失败。这是因为表的辅助索引中的数据可能还在Insert Buffer中，也就是共享表空间中，所以通过ibd文件进行恢复后，还需要进行REPAIR TABLE操作来重建表上所有的辅助索引。

Insert Buffer是一棵B+树，因此其也由叶节点和非叶节点组成。非叶节点存放的是查询的search key（键值），其构造如图2-3所示。



| space | marker | offset |

图 2-3　Insert Buffer非叶节点中的search key

search key一共占用9字节，其中space表示待插入记录所在表的表空间id，在InnoDB存储引擎中，每个表有一个唯一的space id，可以通过space id查询得到是哪张表。space占用4字节。marker占用1字节，它是用来兼容老版本的Insert Buffer。offset表示页所在的位置，占用4字节。

而查询的就是待查询记录的space、offset的内容（即图中的两个数字），InnoDB存储引擎当将待插入记录的search key构造好后，将其插入到Insert Buffer的B+树中，并将记录插入到Insert Buffer B+树的非叶节点中。

对于非叶节点的Insert Buffer B+树的叶节点记录，如图2-4所示，它比叶节点中多一个字段，记录每个辅助索引页中的记录。



图 2-4  Insert Buffer叶节点中的记录

space、marker、page_no这三个字段的含义和上面讨论的一致，一共占用9字节。第4个字段metadata占用4字节，其存储的内容如表2-2所示。

表 2-2  metadata 字段存储的内容

| 名　称 | 字　节 |
| --- | --- |
| IBUF_REC_OFFSET_COUNT | 2 |
| IBUF_REC_OFFSET_TYPE | 1 |
| IBUF_REC_OFFSET_FLAGS | 1 |

IBUF_REC_OFFSET_COUNT是保存两个字节的整数，用来排序每个记录进入Insert Buffer的顺序。因为从InnoDB1.0.x开始支持Change Buffer，这样能保证记录插入到Insert Buffer后，还能源源不断地进行replay操作，继续插入到表中。

当Insert Buffer中的一个辅助5索引记录需要合并到真正的辅助索引中时，需要先进行排序（这也是Insert Buffer B+树在插入记录时，排序的第13个字段的意思）。

最后，在讨论Insert Buffer的插入过程时，我们讨论了space、page_no，它们分别表示待插入记录所在的Insert Buffer B+树所处的位置。因此在进行Merge Insert Buffer操作时，还需要知道另一个重要信息，即表空间中对应的space、page_no的页中哪些页中有记录，这需要用到Insert Buffer Bitmap。

每个Insert Buffer Bitmap用来追踪16384个辅助索引页，也就是256个区（Extent）。每个Insert Buffer Bitmap都在16384个页的第二个页中。每个Insert Buffer Bitmap记录了辅助索引页的相关信息。

每个辅助索引页在Insert Buffer Bitmap中占用4位的bit，其中表2-3所显示了其含义。

表2-3　每个辅助索引页在 Insert Buffer Bitmap 中存储的信息

| 名　称 | 大小（bit） | 说　明 |
|---|---|---|
| IBUF_BITMAP_FREE | 2 | 表示该辅助索引页中的可用空间数量，可取值为：<br>❏ 0 表示无可用剩余空间<br>❏ 1 表示剩余空间大于 1/32 页（512 字节）<br>❏ 2 表示剩余空间大于 1/16 页<br>❏ 3 表示剩余空间大于 1/8 页 |
| IBUF_BITMAP_BUFFERED | 1 | 1 表示该辅助索引页有记录被缓存在 Insert Buffer B+ 树中 |
| IBUF_BITMAP_IBUF | 1 | 1 表示该页为 Insert Buffer B+ 树的索引页 |

## 4.Merge Insert Buffer

通过前面的小节读者应该已经知道了Insert/Change Buffer是一棵B+树。若需要实现插入记录的辅助索引页不在缓冲池中，那么需要将辅助索引记录首先插入到这棵B+树中。但是Insert Buffer中的记录何时合并（merge）到真正的辅助索引中呢？这仍然是我们感兴趣的问题。

概括来说，Merge Insert Buffer的操作可能发生在以下几种情况下：

❏辅助索引页被读取到缓冲池时；

❏Insert Buffer Bitmap页追踪到该辅助索引页已无可用空间时；

❏Master Thread。

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SELECT□□□□□□□□□□□□□□
Insert Buffer Bitmap□□□□□□□□□□□□□□□□□□□□□□□□□□Insert Buffer B+□
□□□□□□□□□Insert Buffer B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

Insert Buffer Bitmap□□□□□□□□□□□□□□□□□□□□□□□□□□□□1/32□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□1/32□□□□□□□□□□□□□□□□□□□□□□□□
□□□□Insert Buffer B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□

□□□□□□□□□□□□□□□Master Thread□□□□□□□Master Thread□□□□□□□□10
□□□□□□□Merge Insert Buffer□□□□□□□□□□□□□□□□merge□□□□□□□□□□□
□□

□Master Thread□□□□□merge□□□□□□□□□□□□□□□□□
srv_innodb_io_capactiy□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□Insert Buffer B+□□□□□□□□□□□□□space□offset□□□□□□□□□□□□□□□□
□space□offset□□□□□□□□□□□□□□□□□□□□□□□Insert Buffer□□□□□□InnoDB
□□□□□□□□□□□□□□□□□□□□□□□□Insert Buffer B+□□□□□□□□□□□□□□□space□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□merge□□□□□□merge□□
□□□□□□□□□□□□□□□□□□□□□Insert/Change Buffer□□□□□□□

## 2.6.2 □□□□

□□□Insert Buffer□□□InnoDB□□□□□□□□□□□□□□□□□□doublewrite□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□16KB□□□□□□□□□□4KB□□□□□□□□□□□□□□□□□□□□□□□□□□□□partial page write□□□□□InnoDB□□□□□□□□doublewrite□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□800□□□'aaaa'□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□apply□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□doublewrite□□□InnoDB□□□□□□□doublewrite□□□□□□□□2-5□□□□

图 2-5　InnoDB体系结构doublewrite工作

doublewrite由两部分组成，一部分是内存中的doublewrite buffer，大小为2MB，另一部分是物理磁盘上共享表空间中连续的128个页，即2个区（extent），大小同样为2MB。在对缓冲池的脏页进行刷新时，并不直接写磁盘，而是会通过memcpy函数将脏页先复制到内存中的doublewrite buffer，之后通过doublewrite buffer再分两次，每次1MB顺序地写入共享表空间的物理磁盘上，然后马上调用fsync函数，同步磁盘，避免缓冲写带来的问题。在这个过程中，因为doublewrite页是连续的，因此这个过程是顺序写的，开销并不是很大。在完成doublewrite buffer页的写入后，再将doublewrite buffer中的页写入各个表空间文件中，此时的写入则是离散的。可以通过以下命令观察到doublewrite运行的情况：

---

mysql∏SHOW GLOBAL STATUS LIKE'innodb_dblwr%'\G;

***************************1.row***************************

Variable_name:Innodb_dblwr_pages_written

Value:6325194

***************************2.row***************************

Variable_name:Innodb_dblwr_writes

Value:100399

2 rows in set(0.00 sec)

---

可以看到，doublewrite一共写了6 325 194个页，但实际的写入次数为100 399，基本上符合64:1。如果发现你的系统在高峰时Innodb_dblwr_pages_written:Innodb_dblwr_writes远小于64:1，那么说明你的系统写入压力并不是很高。

如果操作系统在将页写入磁盘的过程中发生了崩溃，在恢复过程中，InnoDB存储引擎可以从共享表空间中的doublewrite中找到该页的一个副本，将其复制到表空间文件，再应用重做日志。下面显示了一个由doublewrite进行恢复的情况：

---

090924 11:36:32 mysqld restarted

090924 11:36:33 InnoDB:Database was not shut down normally!

InnoDB:Starting crash recovery.

InnoDB:Reading tablespace information from the.ibd files...

InnoDB:Crash recovery may have failed for some.ibd files!

```
InnoDB:Restoring possible half-written data pages from the doublewrite

InnoDB:buffer...
```

---

通过对MySQL数据库的监控（执行SHOW GLOBAL STATUS）Innodb_buffer_pool_pages_flushed，可以观察到有多少页从缓冲池刷新到磁盘中，然而这个参数并不能代表doublewrite的页数量，正确的是Innodb_dblwr_pages_written。在早期的MySQL 5.5.24版本前，由于Innodb_buffer_pool_pages_flushed总是Innodb_dblwr_pages_written的2倍，这是Bug，在MySQL5.5.24进行修复，因此统计这个页数一定要注意，不要被刷新到磁盘的页数Innodb_dblwr_pages_written，具体可以参见MySQL官方手册关于此的描述。

参数skip_innodb_doublewrite可以禁止doublewrite功能，此时可能会发生前面提及的写失效问题。不过如果你有多个slave server，需要提供较快的性能，如slaves erver上的RAID0，也许启用这个参数是一个办法。不过对于需要数据高可靠性的master server，任何时候都应确保开启doublewrite功能。

注意 有些文件系统本身就提供了部分写失效的防范机制，如ZFS文件系统。在这种情况下，用户就不要启用doublewrite了。

# 2.6.3　哈希算法簡介

由於對hash后的主鍵的搜索效率很高，所以在表數據量很大時，以O(1)的時間複雜度，比當前主流索引結構的B+樹有著更高的效率，但是B+樹更加穩定，比如主鍵的B+樹層高一般為3、4層，所以也是3、4次磁盤IO。

InnoDB存儲引擎除了支持我們常用的索引外，還支持一種特殊的索引類型，叫作自適應哈希索引（Adaptive Hash Index，AHI）。AHI是通過緩衝池的B+樹頁構造而來的，所以建立的速度很快，而且也不需要對整張表建立哈希索引。InnoDB存儲引擎會自動根據訪問的頻率和模式來自動地為某些頁建立哈希索引。

AHI有一個要求，即對這個頁的連續訪問模式必須是一樣的。例如對於（a，b）訪問模式情況，其中訪問模式是指查詢的條件一樣。

□WHERE a=xxx

□WHERE a=xxx and b=xxx

若交替進行上述兩種查詢，那麼InnoDB存儲引擎是不會對該頁構造InonDB存儲引擎會自動根據訪問的頻率和模式來自動地為某些頁建立AHI。其中AHI建立的要求如下所示。

□對這個頁連續訪問100次。

□頁通過該模式訪問了N次，其中N=頁中記錄*1/16

由於InnoDB存儲引擎建立的哈希索引，使得AHI對字典類型的查找的速度提升了2倍，對於聯接操作的速度也提高了5倍。毫無疑問，AHI是非常好的優化模式，其設計思想是數據庫自優化（self-tuning），即無須DBA對數據庫進行人為調整。

我們可以通過SHOW ENGINE INNODB STATUS命令來查看當前AHI的使用情況。

---

```
mysql□SHOW ENGINE INNODB STATUS\G;

***************************1.row***************************

Status:

==================================
```

```
090922 11:52:51 INNODB MONITOR OUTPUT

====================================

Per second averages calculated from the last 15 seconds

……

-------------------------------------

INSERT BUFFER AND ADAPTIVE HASH INDEX

-------------------------------------

Ibuf:size 2249,free list len 3346,seg size 5596,

374650 inserts,51897 merged recs,14300 merges

Hash table size 4980499,node heap has 1246 buffer(s)

1640.60 hash searches/s,3709.46 non-hash searches/s

……
```

从上面的例子中可以看到AHI的使用信息，启用AHI后，读取和写入速度可以提高。这里还能看到AHI的使用状态信息，特别需要注意的是哈希索引只能用来搜索等值的查询，如SELECT*FROM table WHERE index_col='xxx'。而对于其他查找类型，如范围查找，是不能使用哈希索引的，因此这里会出现non-hash searches/s的情况。通过hash searches:non-hash searches可以大概了解使用哈希索引后的效率。

由于AHI是由InnoDB存储引擎自动控制的系统管理，但我们可以通过观察，如上面通过SHOW ENGINE INNODB STATUS的结果及参数innodb_adaptive_hash_index来考虑是否需要禁用或启用此特性，默认AHI为开启状态。

# 2.6.4　□□IO

□□□□□□□□□□□□□□□□□□□□□□□□□□IO（Asynchronous IO，AIO）□□□□□□□□□□□□□InnoDB□□□□□□□□□□□

□AIO□□□□□Sync IO□□□□□□□IO□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□IO□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□IO□□□□□□□□□□□□□□□IO□□□□□□IO□□□□□□□□□□□IO□□□□□□□□□AIO□

AIO□□□□□□□□□□□□IO Merge□□□□□□□□□□□IO□□□□1□IO□□□□□□□□IOPS□□□□□□□□□□□□□□□□□□space□page_no□□□□

□8□6□□□□8□7□□□□8□8□

□□□□□□□□□16KB□□□□□□IO□□□□□3□IO□□□□AIO□□□□□□□□□□□□□□□□□□□□□□□□space□page_no□□□□□□□□AIO□□□□□□□□IO□□□□□8□6□□□□□□□48KB□□□□

□□□□Linux□□□□□□□iostat□□□□□□□□□□rrqm/s□wrqm/s□□□□□

---

```
avg-cpu:%user%nice%system%iowait%steal%idle

4.70 0.00 1.60 13.20 0.00 80.50

Device:rrqm/s wrqm/s r/s w/s rMB/s wMB/s avgrq-sz avgqu-sz await svctm%util

sdc 3905.67 172.00 6910.33 466.67 168.81 18.15 51.91 19.17 2.59 0.13 97.73
```

---

□InnoDB1.1.x□□□□AIO□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB 1.1.x□□□□InnoDB Plugin□□□□□□□□□□□□□□□□AIO□□□□□□□Native AIO□□□□□□□□□□□□□□□□□□MySQL□□□□□libaio□□□□□□□□□□□□□□□□□□□□□□□

---

```
/usr/local/mysql/bin/mysqld:error while loading shared libraries:libaio.so.1:cannot open shared object file:No such file or directory
```

---

需要注意的是，Native AIO需要操作系统提供支持。Windows系统和Linux系统都提供Native AIO支持，而Mac OSX系统则未提供。因此在这些系统下，用户可以选择需要针对每个系统进行测试。参数innodb_use_native_aio用来控制是否启用Native AIO，在Linux操作系统下默认值为ON。

---

mysql＞SHOW VARIABLES LIKE'innodb_use_native_aio'\G;

****************************1.row***************************

Variable_name:innodb_use_native_aio

Value:ON

1 row in set(0.00 sec)

---

用户可以通过开启和关闭Native AIO功能来比较InnoDB性能的提升。官方的测试显示，启用Native AIO，恢复速度可以提高75%。

在InnoDB存储引擎中，read ahead方式的读取都是通过AIO完成，脏页的刷新，即磁盘的写入操作则全部由AIO完成。

## 2.6.5   □□□□□□

InnoDB□□□□□□□□□□Flush Neighbor Page□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□extent□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□AIO□□□□□□IO□□□□□□□□□□□IO□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□IOPS□□□□□□□□□□□□□□□□

□□□□InnoDB□□□□□□1.2.x□□□□□□□□□□□innodb_flush_neighbors□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□IOPS□□□□□□□□□□□□□□□□□□0□□□□□□□□□□□□

## 2.7 □□□□□□□□□

InnoDB□MySQL□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□
MySQL□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□

□□□□□□□□□innodb_fast_shutdown□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□
□□□0□1□2□□□□□□1□

□0□□□□MySQL□□□□□□□□□InnoDB□□□□□□□□□full purge□merge
insert buffer□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□InnoDB□□□□□□□□□□□□□□□0□□□□□□□□□□□□□

□1□□□□innodb_fast_shutdown□□□□□□□□□□□□□□□□□□□□full purge□
merge insert buffer□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□2□□□□□□full purge□merge insert buffer□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□
□□□recovery□□□

□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□"□□"□□□□□□□□□□□□□□□□□□□□□kill□
□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
innodb_fast_shutdown□□□□□2□□□□□□MySQL□□□□□□□□□□□□□InnoDB□□□□□
□□□□□□□□□□□□□□

□□□innodb_force_recovery□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□
0□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□corruption□
MySQL□□□□□□□□□□□□□crash□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
alter table□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□

□□□innodb_force_recovery□□□□□□□□6□□□□□□1□6□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□

□1(SRV_FORCE_IGNORE_CORRUPT)：忽略检查到的corrupt页。

□2(SRV_FORCE_NO_BACKGROUND)：阻止Master Thread线程的运行，如Master Thread线程需要进行full purge操作，而这会导致crash。

□3(SRV_FORCE_NO_TRX_UNDO)：不进行事务的回滚操作。

□4(SRV_FORCE_NO_IBUF_MERGE)：不进行插入缓冲的合并操作。

□5(SRV_FORCE_NO_UNDO_LOG_SCAN)：不查看撤销日志（Undo Log），InnoDB存储引擎会将未提交的事务视为已提交。

□6(SRV_FORCE_NO_LOG_REDO)：不进行前滚的操作。

需要特别注意的是，在设置参数innodb_force_recovery大于0后，用户可以对表进行select、create、drop操作，但insert、update和delete这类DML操作是不允许的。

因此，若出现下面提示信息，即用户正在进行的session需要对1 000万行的InnoDB表进行更新操作，但是该操作显然不会成功：

---

mysql＞START TRANSACTION;

Query OK,0 rows affected(0.00 sec)

mysql＞UPDATE Profile SET password='';

Query OK,9587770 rows affected(7 min 55.73 sec)

Rows matched:9999248 Changed:9587770 Warnings:0

---

START TRANSACTION语句开始了一个事务，因为事务不自动提交（auto commit），所以UPDATE操作产生了大量的UNDO。由于undo log变得很大，用户选择用kill操作结束MySQL服务进程，如下：

---

[root@nineyou0-43～]#ps-ef|grep mysqld

root 28007 1 0 13:40 pts/1 00:00:00/bin/sh./bin/mysqld_safe--datadir=/usr/local/mysql/data--pid-file=/usr/local/mysql/data/nineyou0-43.pid

mysql 28045 28007 42 13:40 pts/1 00:04:23/usr/local/mysql/bin/mysqld--basedir=/usr/local/mysql--datadir=/usr/local/mysql/data--user=mysql--pid-file=/usr/local/mysql/data/nineyou0-43.pid--skip-external-locking--port=3306--socket=/tmp/mysql.sock

root 28110 26963 0 13:50 pts/11 00:00:00 grep mysqld

[root@nineyou0-43 ]#kill-9 28007

[root@nineyou0-43 ]#kill-9 28045

---

## 通过kill命令关闭之后我们再次启动MySQL数据库。由于之前的UPDATE操作因为实例被强制终止了，这时肯定需要回滚操作，而这可以通过查看err日志文件来观察，如下显示了这个时候的日志：

---

090922 13:40:20 InnoDB:Started;log sequence number 6 2530474615

InnoDB:Starting in background the rollback of uncommitted transactions

090922 13:40:20 InnoDB:Rolling back trx with id 0 5281035,8867280 rows to undo

InnoDB:Progress in percents:1090922 13:40:20

090922 13:40:20[Note]/usr/local/mysql/bin/mysqld:ready for connections.

Version:'5.0.45-log'socket:'/tmp/mysql.sock'port:3306 MySQL Community Server(GPL)

2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100

InnoDB:Rolling back of trx id 0 5281035 completed

090922 13:49:21 InnoDB:Rollback of non-prepared transactions completed

---

## 从日志中可以看到，由于这时参数innodb_force_recovery设为0，InnoDB会把所有未提交的事务全部回滚，而这时共需要回滚8867280行记录，整个回滚操作需要9分钟。

## 接着我们来做第二种情况的测试，我们重启MySQL数据库，并把参数innodb_force_recovery设为3，然后再次将InnoDB存储引擎进行回滚操作，在错误日志文件中会看到如下内容：

---

090922 14:26:23 InnoDB:Started;log sequence number 7 2253251193

InnoDB:!!!innodb_force_recovery is set to 3!!!

090922 14:26:23[Note]/usr/local/mysql/bin/mysqld:ready for connections.

Version:'5.0.45-log'socket:'/tmp/mysql.sock'port:3306 MySQL Community Server(GPL)

---

□□□□□"!!!"，InnoDB□□□□□innodb_force_recovery□□□□3□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 2.8 □□

□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□"□□"□□□□□□□□□□□□□MySQL□□□□□□□□□□□□InnoDB□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□3□□□□□□MySQL□□□□□□□MySQL□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□

# □3□　□□

□□□□□□MySQL□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□

□socket□□□□□□UNIX□□□□□□□□□□□□□□□□□□□□□□□

□pid□□□□MySQL□□□□□□ID□□□□

□MySQL□□□□□□□□□□□MySQL□□□□□□□□□□

□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□

# 3.1 参数概述

从1.4节的示例中可以看到，MySQL实例在启动时会去读取配置文件中的参数。如果没有配置文件，那么启动时实例会按照默认的参数设置启动。在下面的例子中，尝试使用命令mysql--help|grep my.cnf来寻找配置文件。

MySQL数据库参数文件的作用和Oracle数据库的参数文件极其类似，不同的是，Oracle实例在启动时若找不到参数文件，是不能进行实例的mount操作的，而MySQL数据库却可以，MySQL可以不需要参数文件，这时所有的参数值取决于编译MySQL时指定的默认值和源代码中指定参数的默认值。但是，如果MySQL实例在默认的数据库目录下找不到mysql架构，则启动同样会失败，此时可能在错误文件中找到如下内容：

---

```
090922 16:25:52 mysqld started

090922 16:25:53 InnoDB:Started;log sequence number 8 2801063211

InnoDB:!!!innodb_force_recovery is set to 1!!!

090922 16:25:53[ERROR]Fatal error:Can't open and lock privilege tables:Table'mysql.host'doesn't exist

090922 16:25:53 mysqld ended
```

---

MySQL中mysql架构中记录了访问该实例的权限，当找不到这个架构时，MySQL实例不能成功启动。

MySQL数据库的参数文件是以文本方式进行存储的。可以直接通过一些常用的文本编辑软件（如vi和emacs）进行参数的修改。

## 3.1.1 什么是参数

简单地说，可以把数据库参数看成一个键/值对（key/value）。1.4节已经显示了一个MySQL的InnoDB存储引擎内存的变量innodb_buffer_pool_size，我们可以看到它被设置的值为1G，即innodb_buffer_pool_size=1G，其中"键"为innodb_buffer_pool_size，"值"为1G。这就是一个键值对。可以通过SHOW VARIABLES查看数据库中的所有参数，也可以通过LIKE来过滤参数名。从MySQL 5.1版本开

# □□□□□□□□information_schema□□□□□GLOBAL_VARIABLES□□□□□□□□□□□□□□□

---

```
mysql□SELECT*FROM

-□GLOBAL_VARIABLES

-□WHERE VARIABLE_NAME LIKE'innodb_buffer%'\G;

***************************1.row***************************

VARIABLE_NAME:INNODB_BUFFER_POOL_SIZE

VARIABLE_VALUE:1073741824

1 row in set(0.00 sec)

mysql□SHOW VARIABLES LIKE'innodb_buffer%'\G;

***************************1.row***************************

Variable_name:innodb_buffer_pool_size

Value:1073741824

1 row in set(0.00 sec)
```

---

## □□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□GLOBAL_VARIABLES□□□□□□□□□□□□□□□□□□□SHOW VARIABLES□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□

Oracle□□□□□□□□□□□□□□□undocumented parameter□□□□□Oracle"□□□□"□□□SQL Server□□□□□□□□□□□DBA□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Oracle□SQL Server□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 3.1.2 　系统变量

MySQL中的系统变量分为以下两种类型：

□动态（dynamic）变量。

□静态（static）变量。

动态变量意味着在MySQL服务实例运行期间它们的值可以通过命令行进行修改；静态变量意味着在整个实例的生命周期内它是只读（read only）的，可以用命令SET修改系统变量的值。修改语法：SET命令格式如下：

---

SET

|[global|session]system_var_name=expr

|[@@global.|@@session.|@@]system_var_name=expr

---

可以使用选项global、session为系统变量分别设置全局级别和会话级别的值。一些系统变量只有全局值，如变量autocommit；一些系统变量只有会话值，如变量binlog_cache_size；一些系统变量既拥有全局值又拥有会话值，如变量read_buffer_size。举例如下：

---

mysql>SET read_buffer_size=524288;

Query OK,0 rows affected(0.00 sec)

mysql>SELECT@@session.read_buffer_size\G;

***************************1.row***************************

@@session.read_buffer_size:524288

1 row in set(0.00 sec)

mysql>SELECT@@global.read_buffer_size\G;

***************************1.row***************************

@@global.read_buffer_size:2093056

1 row in set(0.00 sec)

---

□□□□□□□□□□□□□read_buffer_size由2MB修改为512KB，而在全局级别上的read_buffer_size仍然为2MB。如果会话修改完成之后，重新登录MySQL数据库，read_buffer_size由原2MB恢复至512KB。因此，使用set global|session命令修改系统变量值，实际上等同于使用SET@@globl|@@session命令修改系统变量

---

mysql□SET@@global.read_buffer_size=1048576;

Query OK,0 rows affected(0.00 sec)

mysql□SELECT@@session.read_buffer_size\G;

***************************1.row***************************

@@session.read_buffer_size:524288

1 row in set(0.00 sec)

mysql□SELECT@@global.read_buffer_size\G;

***************************1.row***************************

@@global.read_buffer_size:1048576

1 row in set(0.00 sec)

---

会话级read_buffer_size的大小仍然为1MB，而全局级read_buffer_size则被修改512KB。需要说明的是，并不是所有的系统变量均可以通过这种方式进行设置，有些系统变量只能是全局的，有些系统变量可以存在于全局范围和会话范围，而有些系统变量既不属于全局范围又不属于会话范围。另外，并不是所有的系统变量都可以进行动态修改，有关MySQL系统变量的详细信息，请参见MySQL官方文档Dynamic System Variables部分的内容。

如下所示的系统变量只能读取，不能设置。

---

mysql□SET GLOBAL datadir='/db/mysql';

ERROR 1238(HY000):Variable'datadir'is a read only variable

---

## 3.2 　日志管理

本章会重点介绍与MySQL运维管理密切相关的几种MySQL日志，它们分别是下面四种。

□ 错误日志（error log）

□ 二进制日志（binlog）

□ 慢查询日志（slow query log）

□ 查询日志（log）

了解这些日志，对于DBA对MySQL的日常管理非常有用，本节会逐一介绍这几种日志。

## 3.2.1 　错误日志

错误日志记录MySQL数据库系统的诊断信息，这对于MySQL DBA进行故障诊断和排错非常有用。系统出现任何错误信息，都会记录到错误日志中。可以通过命令查看数据库的错误日志（SHOW VARIABLES LIKE'log_error'），结果如下所示：

---

```
mysql＞SHOW VARIABLES LIKE'log_error'\G;

***************************1.row***************************

Variable_name:log_error

Value:/mysql_data_2/stargazer.log

1 row in set(0.00 sec)

mysql＞system hostname

stargazer
```

---

从上面的例子可以看出，错误日志的文件名是由参数指定的。在默认情况下，错误日志文件名为主机名（stargazer），后缀名为默认的startgazer.err。当数据库MySQL发生故障无法正常使用时，可以首先查看此文件，以便进行错误定位。此文件不仅记录所有的错误信息，也记录了一些警告信息或者正确的信息。

---

[root@nineyou0-43 data]#tail-n 50 nineyou0-43.err

090924 11:31:18 mysqld started

090924 11:31:18 InnoDB:Started;log sequence number 8 2801063331

090924 11:31:19[ERROR]Fatal error:Can't open and lock privilege tables:Table'mysql.host'doesn't exist

090924 11:31:19 mysqld ended

从上面的例子看出，由于使用了mysql数据库中不存在的存储引擎，导致启动失败，并且出现了非常明确的错误信息。这里的warning信息就是前面提到的，因为存储引擎的更换，忽略了原表的创建参数，重新以默认的存储引擎InnoDB重建了redo log。

090924 11:39:44 InnoDB:ERROR:the age of the last checkpoint is 9433712,

InnoDB:which exceeds the log group capacity 9433498.

InnoDB:If you are using big BLOB or TEXT rows,you must set the

InnoDB:combined size of log files at least 10 times bigger than the

InnoDB:largest such row.

090924 11:40:00 InnoDB:ERROR:the age of the last checkpoint is 9433823,

InnoDB:which exceeds the log group capacity 9433498.

InnoDB:If you are using big BLOB or TEXT rows,you must set the

InnoDB:combined size of log files at least 10 times bigger than the

InnoDB:largest such row.

090924 11:40:16 InnoDB:ERROR:the age of the last checkpoint is 9433645,

InnoDB:which exceeds the log group capacity 9433498.

InnoDB:If you are using big BLOB or TEXT rows,you must set the

InnoDB:combined size of log files at least 10 times bigger than the

InnoDB:largest such row.

# 3.2.2　慢查询日志

3.2.1小节介绍了错误日志，本小节将介绍第二个日志，即慢查询日志（slow log）。慢查询日志可帮助DBA定位可能存在问题的SQL语句，从而进行SQL语句层面的优化。例如，可以在MySQL启动时设一个阈值，将运行时间超过该值的所有SQL语句都记录到慢查询日志文件中。DBA每天或每过一段时间对运行日志所记录的SQL语句进行分析，看是否有SQL语句需要进行优化。该阈值可以通过参数long_query_time来设置，默认值为10，代表10秒。

在默认情况下，MySQL数据库并不启动慢查询日志，用户需要手工将这个参数设为ON：

```
mysql：SHOW VARIABLES LIKE'long_query_time'\G;

***************************1.row***************************

Variable_name:long_query_time

Value:10.000000

1 row in set(0.00 sec)

mysql：SHOW VARIABLES LIKE'log_slow_queries'\G;

***************************1.row***************************

Variable_name:log_slow_queries

Value:ON

1 row in set(0.00 sec)
```

这里有一点需要注意，设置long_query_time这个阈值后，MySQL数据库会记录运行时间超过该值的所有SQL语句，但运行时间正好等于long_query_time的情况并不会被记录下。也就是说，在源代码里是判断大于long_query_time，而非大于等于。从MySQL 5.1开始，long_query_time开始以微秒记录SQL语句运行的时间，之前仅用秒为单位记录。而这样可以更精确地记录SQL的运行时间，供DBA分析。对DBA来说，一条SQL语句运行0.5秒和0.05秒的区别可能是非常大的，因此有时候可能需要微秒级别的记录。

另一个和慢查询日志有关的参数是log_queries_not_using_indexes，如果运行的SQL语句没有使用索引，则MySQL数据库同样会将这条SQL语句记录到慢查询日志文件。首先确认打开了log_queries_not_using_indexes：

```
mysql□SHOW VARIABLES LIKE'log_queries_not_using_indexes'\G;

***************************1.row***************************

Variable_name:log_queries_not_using_indexes

Value:ON

1 row in set(0.00 sec)
```

---

# MySQL 5.6.5版本以后还引入了一个参数log_throttle_queries_not_using_indexes。该参数用来表示每分钟允许记录到slow log的且未使用索引的SQL语句次数。该值默认为0，表示没有限制。在生产环境下，若没有使用索引，此类SQL语句会频繁地被记录到slow log，从而导致slow log文件的大小不断增加，故DBA可通过此参数进行配置。

# DBA可以像一般的查询语句那样分析慢查询SQL语句。但是更简单直接的方法是使用MySQL数据库提供的慢查询日志分析工具，它能帮助用户更快速地定位慢查询SQL语句，从而为优化赢得时间。在本书中，我们介绍MySQL数据库提供的mysqldumpslow工具，这也是最常用的DBA分析工具。

---

```
[root@nh122-190 data]#mysqldumpslow nh122-190-slow.log

Reading mysql slow query log from nh122-190-slow.log

Count:11 Time=10.00s(110s)Lock=0.00s(0s)Rows=0.0(0),dbother[dbother]@localhost

insert into test.DbStatus select now(),(N-com_select)/(N-uptime),(N-com_insert)/(N-uptime),(N-com_update)/(N-uptime),(N-com_delete)/(N-uptime),N-(N/N),N-(N/N),N.N/N,N-N/(N*N),GetCPULoadInfo(N)from test.CheckDbStatus order by check_id desc limit N

Count:653 Time=0.00s(0s)Lock=0.00s(0s)Rows=0.0(0),9YOUgs_SC[9YOUgs_SC]@[192.168.43.7]

select custom_name_one from'low_game_schema'.'role_details'where role_id='S'rse and summarize the MySQL slow query log.Options are

--verbose verbose

--debug debug

--help write this text to standard output

-v verbose

-d debug

-s ORDER what to sort by(al,at,ar,c,l,r,t),'at'is default

al:average lock time

ar:average rows sent
```

at:average query time

c:count

l:lock time

r:rows sent

t:query time

-r reverse the sort order(largest last instead of first)

-t NUM just show the top n queries

-a don't abstract all numbers to N and strings to'S'

-n NUM abstract numbers with at least n digits within names

-g PATTERN grep:only consider stmts that include this string

-h HOSTNAME hostname of db server for*-slow.log filename(can be wildcard),

default is'*',i.e.match all

-i NAME name of server instance(if using mysql.server startup script)

-l don't subtract lock time from total time

---

# 下面是工作常用的一个命令，取出10条SQL语句，同时按照时间排序：

---

[root@nh119-141 data]#mysqldumpslow-s al-n 10 david.log

Reading mysql slow query log from david.log

Count:5 Time=0.00s(0s)Lock=0.20s(1s)Rows=4.4(22),Audition[Audition]@[192.168.30.108]

SELECT OtherSN,State FROM wait_friend_info WHERE UserSN=N

Count:1 Time=0.00s(0s)Lock=0.00s(0s)Rows=1.0(1),audition-kr[audition-kr]@[192.168.30.105]

SELECT COUNT(N)FROM famverifycode WHERE UserSN=N AND verifycode='S'

……

---

# MySQL 5.1版本以后，慢查询日志可以写到数据库中，这样我们可以更方便地查看慢查询日志，在mysql库中有个slow_log的表。表结构如下：

---

mysql｝SHOW CREATE TABLE mysql.slow_log\G;

***************************1.row***************************

Table:slow_log

Create Table:CREATE TABLE'slow_log'(

'start_time'timestamp NOT NULL DEFAULT CURRENT_TIMESTAMP ON UPDATE CURRENT_TIMESTAMP,

'user_host'mediumtext NOT NULL,

'query_time'time NOT NULL,

'lock_time'time NOT NULL,

'rows_sent'int(11)NOT NULL,

'rows_examined'int(11)NOT NULL,

'db'varchar(512)NOT NULL,

'last_insert_id'int(11)NOT NULL,

'insert_id'int(11)NOT NULL,

'server_id'int(11)NOT NULL,

'sql_text'mediumtext NOT NULL

)ENGINE=CSV DEFAULT CHARSET=utf8 COMMENT='Slow log'

1 row in set(0.00 sec)

## 参数log_output指定了慢查询输出的格式，默认为FILE，可以将它设置为TABLE，然后就可以查询mysql架构下的slow_log表了，如下所示：

mysql□SHOW VARIABLES LIKE'log_output'\G;

+---------------+---------+

|Variable_name|Value|

+---------------+---------+

|log_output|FILE|

+---------------+---------+

1 row in set(0.00 sec)

mysql□SET GLOBAL log_output='TABLE';

Query OK,0 rows affected(0.00 sec)

mysql□SHOW VARIABLES LIKE'log_output'\G;

+---------------+---------+

```
|Variable_name|Value|

+---------------+---------+

|log_output|TABLE|

+---------------+---------+

1 row in set(0.00 sec)

mysql□select sleep(10)\G;

+-----------+

|sleep(10)|

+-----------+

|0|

+-----------+

1 row in set(10.01 sec)

mysql□SELECT*FROM mysql.slow_log\G;

***************************1.row***************************

start_time:2009-09-25 13:44:29

user_host:david[david]@localhost[]

query_time:00:00:09

lock_time:00:00:00

rows_sent:1

rows_examined:0

db:mysql

last_insert_id:0

insert_id:0

server_id:0

sql_text:select sleep(10)

1 row in set(0.00 sec)
```

# 当□log_output□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□sleep□10□□□□□□SQL□□□□□□□□□slow_log□□□□

由于slow_log表不允许修改，因此需要将CSV存储引擎的慢查询日志表转换为可以修改的slow_log表，修改为MyISAM，对start_time列添加索引等。不过在进行修改之前需要确认当前没有开启日志记录功能，否则会报错。

---

mysql□ALTER TABLE mysql.slow_log ENGINE=MyISM;

ERROR 1580(HY000):You cannot'ALTER'a log table if logging is enabled

mysql□SET GLOBAL slow_query_log=off;

Query OK,0 rows affected(0.00 sec)

mysql□ALTER TABLE mysql.slow_log ENGINE=MyISAM;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

---

在上述操作过程中，由于slow_log表本身默认是不支持MyISAM存储引擎的，所以只能通过将日志记录功能关闭后，再进行存储引擎类型的修改操作。

MySQL□slow log能够记录执行时间超过指定阈值的所有查询，这对于定位产生性能问题的查询是很有帮助的。但是如果需要对数据库中执行的所有SQL进行统计，找出其中执行时间超过0.5秒的

InnoSQL数据库记录的SQL日志相比，不同之处是，MySQL原始的slow log仅仅记录了逻辑读取logical reads，物理读取physical reads。这里逻辑读取和物理读取的差异就在于读取是否需要进行IO操作。下面显示了一个打印到慢查询日志中的慢查询信息：

---

#Time:111227 23:49:16

#User@Host:root[root]@localhost[127.0.0.1]

#Query_time:6.081214 Lock_time:0.046800 Rows_sent:42 Rows_examined:727558 Logical_reads:91584 Physical_reads:19

use tpcc;

SET timestamp=1325000956;

SELECT orderid,customerid,employeeid,orderdate

FROM orders

WHERE orderdate IN

(SELECT MAX(orderdate)

```
FROM orders

GROUP BY(DATE_FORMAT(orderdate,'%Y%M'))

);
```

---

以上语句一共执行了这么多行，扫描了91 584行数据，使用19秒的时间。通过这个我们就可以定位到慢SQL，进行优化了。

还有一个参数慢查询，叫long_query_io，当查询的扫描IO次数的SQL执行时间在slow log时进行查询。超过100次的扫描，或者执行时间超过100的SQL记录都将在slow log里。我们也可以使MySQL数据库可以记录下这些，这个slow_query_type定义了记录到slow log中的查询类型。

□0：不将SQL语句记录到slow log

□1：根据运行时间将SQL语句记录到slow log

□2：根据逻辑IO次数将SQL语句记录到slow log

□3：根据运行时间以及逻辑IO次数将SQL语句记录到slow log

# 3.2.3 □□□□

□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□.log□□□□□□□□□□□□□□

---

```
[root@nineyou0-43 data]#tail nineyou0-43.log

090925 11:00:24 44 Connect zlm@192.168.0.100 on

44 Query SET AUTOCOMMIT=0

44 Query set autocommit=0

44 Quit

090925 11:02:37 45 Connect Access denied for user'root'@'localhost'(using password:NO)

090925 11:03:51 46 Connect Access denied for user'root'@'localhost'(using password:NO)

090925 11:04:38 23 Query rollback
```

---

□□□□□□□□□□□□□□□□□□□□□□□□□Access denied□□□□□□□□□□□□□□□□□□□□
SQL□□□□□□□□□□□□□□□□□□□□□□□MySQL 5.1□□□□□□□□□□□□□□□□□□□□
mysql□□□□□general_log□□□□□□□□□□□□□□□□□□□□□□□□slow_log□□□□□□□□
□□□□□□□

# 3.2.4 □□□□□

□□□□□□binary log□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□SELECT□SHOW□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□UPDATE t SET a=1 WHERE a=2;

Query OK,0 rows affected(0.00 sec)

Rows matched:0 Changed:0 Warnings:0

mysql□SHOW MASTER STATUS\G;

***************************1.row***************************

File:mysqld.000008

Position:383

Binlog_Do_DB:

Binlog_Ignore_DB:

Executed_Gtid_Set:

1 row in set(0.00 sec)

mysql□SHOW BINLOG EVENTS IN'mysqld.000008'\G;

***************************1.row***************************

Log_name:mysqld.000008

Pos:4

Event_type:Format_desc

Server_id:1

End_log_pos:120

Info:Server ver:5.6.6-m9-log,Binlog ver:4

***************************2.row***************************

Log_name:mysqld.000008

Pos:120

Event_type:Query

Server_id:1

```
End_log_pos:199

Info:BEGIN

***************************3.row***************************

Log_name:mysqld.000008

Pos:199

Event_type:Query

Server_id:1

End_log_pos:303

Info:use'test';UPDATE t SET a=1 WHERE a=2

***************************4.row***************************

Log_name:mysqld.000008

Pos:303

Event_type:Query

Server_id:1

End_log_pos:383

Info:COMMIT

4 rows in set(0.00 sec)
```

从上面的例子中可以看到，MySQL数据库首先进行UPDATE操作，从返回结果看该事务已经Changed了0行，但是通过二进制日志可以明显地发现，在之后的SHOW BINLOG EVENT操作中看到了生成了一个UPDATE的操作。

对于写入性能，SELECT、SHOW这类操作是不会把事务记录到二进制日志，只有更改数据库的操作才会写入到二进制日志。总的来说，二进制日志主要有以下几种作用：

□ 恢复（recovery）：某些数据的恢复需要二进制日志，例如，在一个数据库全备文件恢复后，用户可以通过二进制日志进行point-in-time的恢复。

□ 复制（replication）：其原理与恢复类似，通过复制和执行二进制日志使一台远程的MySQL数据库（一般称为slave或standby）与一台MySQL数据库（一般称为master或primary）进行实时同步。

□ 审计（audit）：用户可以通过二进制日志中的信息来进行审计，判断是否有对数据库

通过使用选项log-bin[=name]可以开启二进制日志，如果不指定name，则默认二进制日志文件名为主机名，后缀名为二进制日志的序列号，所在路径为数据库所在目录（datadir），如：

mysql>show variables like'datadir';

+---------------+---------------------------+

|Variable_name|Value|

+---------------+---------------------------+

|datadir|/usr/local/mysql/data/|

+---------------+---------------------------+

1 row in set(0.00 sec)

mysql>system ls-lh/usr/local/mysql/data/;

total 2.1G

-rw-rw----1 mysql mysql 6.5M Sep 25 15:13 bin_log.000001

-rw-rw----1 mysql mysql 17 Sep 25 00:32 bin_log.index

-rw-rw----1 mysql mysql 300M Sep 25 15:13 ibdata1

-rw-rw----1 mysql mysql 256M Sep 25 15:13 ib_logfile0

-rw-rw----1 mysql mysql 256M Sep 25 15:13 ib_logfile1

drwxr-xr-x 2 mysql mysql 4.0K May 7 10:08 mysql

drwx------2 mysql mysql 4.0K May 7 10:09 test

这里的bin_log.00001即为二进制日志文件，我们在配置文件中指定了名字，所以没有用主机名。bin_log.index为二进制的索引文件，用来存储过去产生的二进制日志序号，在通常情况下，不建议手工修改这个文件。

二进制日志文件在默认情况下并没有启动，需要手动指定参数来启动。可能有人会质疑，开启这个选项是否会对数据库整体性能有所影响。不错，开启这个选项的确会影响性能，但是性能的损失十分有限。根据MySQL官方手册中的测试表明，开启二进制日志会使性能下降1%。但考虑到可以使用复制（replication）和point-in-time的恢复，这些性能损失绝对是可以且应该被接受的。

二进制日志主要有以下几种作用，具体如下所示。

□max_binlog_size

❑binlog_cache_size

❑sync_binlog

❑binlog-do-db

❑binlog-ignore-db

❑log-slave-update

❑binlog_format

□□max_binlog_size□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□+1□□□□□.index□□□□MySQL 5.0□□□□□□□1 073 741 824□□□1 G□□□□□□□max_binlog_size□□□□□1.1G□□□

□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□uncommitted□□□□□□□□□□□□□□□□□□□□□□□□committed□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□binlog_cache_size□□□□□□□□□32K□□□□binlog_cache_size□□□□□□session□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□binlog_cache_size□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□binlog_cache_size□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□SHOW GLOBAL STATUS□□□□binlog_cache_use□binlog_cache_disk_use□□□□□□□□□□□binlog_cache_size□□□□□□□□Binlog_cache_use□□□□□□□□□□□□□□□□binlog_cache_disk_use□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

mysql□show variables like'binlog_cache_size';

+------------------+------+

|Variable_name|Value|

+------------------+------+

|binlog_cache_size|32768|

+------------------+------+

1 row in set(0.00 sec)

```
mysql》show global status like'binlog_cache%';

+----------------------+--------------+

|Variable_name|Value|

+----------------------+--------------+

|binlog_cache_disk_use|0|

|binlog_cache_use|33553|

+----------------------+--------------+

2 rows in set(0.00 sec)
```

□□□□□□□□□□33 553□□□□□□□□□□□□□□0□□□□32KB□□□□□□□□□□□□□□□□MySQL□□□□
□□□□□□□□□□□□□□□□□□binlog_cache_size□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
sync_binlog=[N]□□□□□□□□□□□□□□□□□□□□□□□□□N□□1□□sync_binlog=1
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
sync_binlog□□□□□□0□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□ON□□□□□□□□ON□□□□□□□□□□□□□□IO□□□□□□□□□□□□□

□□□□□□□□sync_binlog□□1□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□
□□□□□□□□□□COMMIT□□□□□□□□□□sync_binlog□1□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□
□□□□□COMMIT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□innodb_support_xa□□1□□□□□□□□□
innodb_support_xa□XA□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□
□□□□□□□

□□□binlog-do-db□binlog-ignore-db□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□slave□□□□□□□□□□□master□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□log-slave-update□□□□□□□□□master=□slave=
□slave□□□□□□□□□□□□□□□□□□□

binlog_format参数十分重要，它影响了记录二进制日志的格式。在MySQL 5.1版本之前，没有这个参数。所有二进制文件的格式都是基于SQL语句（statement）级别的，因此基于这个格式的二进制日志文件的复制（Replication）和Oracle数据库的Standby有些相似。同时，对于复制是有一定要求的。如在主服务器运行rand、uuid等函数，又或者使用触发器等操作，这些都可能会导致主从服务器上表中数据的不一致（not sync）。另一个影响是，会发现InnoDB存储引擎的默认事务隔离级别是REPEATABLE READ。这其实也是因为二进制日志文件格式的关系，如果使用READ COMMITTED的事务隔离级别（大多数数据库，如Oracle、Microsoft SQL Server数据库的默认隔离级别），会出现类似丢失更新的现象，从而出现主从数据库上的数据不一致。

MySQL 5.1开始引入了binlog_format参数，该参数可设的值有STATEMENT、ROW和MIXED。

（1）STATEMENT格式和之前的MySQL版本一样，二进制日志文件记录的是日志的逻辑SQL语句。

（2）在ROW格式下，二进制日志记录的不再是简单的SQL语句了，而是记录表的行更改情况。基于ROW格式的复制类似于Oracle的物理Standby（有些可能认为两者相似）。同时，解决了上述提及的Statement格式下复制的问题。从MySQL 5.1版本开始，如果设置了binlog_format为ROW，可以将InnoDB的事务隔离级别设为READ COMMITTED，以获得更好的并发性。

（3）在MIXED格式下，MySQL默认采用STATEMENT格式进行二进制日志文件的记录，但是在一些情况下会使用ROW格式，可能的情况有：

1．表的存储引擎为NDB，这时对表的DML操作都会以ROW格式记录。

2．使用了UUID()、USER()、CURRENT_USER()、FOUND_ROWS()、ROW_COUNT()等不确定函数。

3．使用了INSERT DELAY语句。

4．使用了用户定义函数（UDF）。

5．使用了临时表（temporary table）。

此外，binlog_format参数还有一些其他约束，如表3-1所示。

表 3-1 存储引擎对二进制日志格式的支持情况

| 存储引擎 | Row 格式 | Statement 格式 |
|---|---|---|
| InnoDB | Yes | Yes |
| MyISAM | Yes | Yes |
| HEAP | Yes | Yes |
| MERGE | Yes | Yes |
| NDB | Yes | No |
| Archive | Yes | Yes |
| CSV | Yes | Yes |
| Federate | Yes | Yes |
| Blockhole | No | Yes |

binlog_format是动态参数,因此可以在数据库运行环境下进行更改,例如,我们可以将当前会话的binlog_format设为ROW,如:

```
mysql>SET@@session.binlog_format='ROW';

Query OK,0 rows affected(0.00 sec)

mysql>SELECT@@session.binlog_format;
```

```
+------------------------------+

|@@session.binlog_format|

+------------------------------+

|ROW|

+------------------------------+

1 row in set(0.00 sec)
```

上面的语句可以看到binlog_format设置的是会话级别的变量，只对当前会话有效。如果要对全局设置可以用下面的语句。

```
mysql：SET GLOBAL binlog_format='ROW';

Query OK,0 rows affected(0.00 sec)

mysql：SELECT@@global.binlog_format;

+------------------------------+

|@@global.binlog_format|

+------------------------------+

|ROW|

+------------------------------+

1 row in set(0.00 sec)
```

我们用全局的语句设置的binlog_format格式是ROW，为什么在会话级别的变量依然显示的是原来的值呢？原因是全局的参数对当前会话是不生效的，ROW格式记录的是每行数据的变化，这样如果是批量操作（比如100W行数据的UPDATE），会生成大量的日志，我们可以看下面的例子。

```
mysql：SELECT@@session.binlog_format\G;

***************************1.row***************************

@@session.binlog_format:STATEMENT

1 row in set(0.00 sec)

mysql：SHOW MASTER STATUS\G;

***************************1.row***************************
```

File:test.000003

Position:106

Binlog_Do_DB:

Binlog_Ignore_DB:

1 row in set(0.00 sec)

mysql␤UPDATE t1 SET username=UPPER(username);

Query OK,89279 rows affected(1.83 sec)

Rows matched:100000 Changed:89279 Warnings:0

mysql␤SHOW MASTER STATUS\G;

***************************1.row***************************

File:test.000003

Position:306

Binlog_Do_DB:

Binlog_Ignore_DB:

1 row in set(0.00 sec)

---

# 可以看到，在binlog_format设置为STATEMENT时，执行这条UPDATE语句所产生的日志量是200字节（306-106）。下面更改为ROW，然后在表t2上执行同样的操作：

---

mysql␤SET SESSION binlog_format='ROW';

Query OK,0 rows affected(0.00 sec)

mysql␤SHOW MASTER STATUS\G;

***************************1.row***************************

File:test.000003

Position:306

Binlog_Do_DB:

Binlog_Ignore_DB:

1 row in set(0.00 sec)

mysql␤UPDATE t2 SET username=UPPER(username);

Query OK,89279 rows affected(2.42 sec)

Rows matched:100000 Changed:89279 Warnings:0

mysql□SHOW MASTER STATUS\G;

***************************1.row***************************

File:test.000003

Position:13782400

Binlog_Do_DB:

Binlog_Ignore_DB:

1 row in set(0.00 sec)

---

可以看到，现在的二进制日志（ROW格式）增加了约13 782 094字节。而之前我们已经对表进行了更新，约13MB的大小，因此t2表的大小变为了17MB。可见二进制日志记录了对每行记录的修改。由于sync_binlog=1的关系，在我们重启MySQL数据库后，日志文件SQL语句依然存在，并没有丢失，也没有发生损坏。

由于二进制日志文件在默认情况下binlog_format为ROW格式，我们可以将其更改为STATEMENT格式来查看之前所运行的更新操作究竟执行了哪些逻辑操作。

二进制日志文件是不能像错误日志文件、慢查询日志文件那样用cat、head、tail等命令来查看的。要查看二进制日志文件的内容，须通过MySQL提供的工具mysqlbinlog。对于STATEMENT格式的二进制日志文件，在使用mysqlbinlog后，看到的就是执行的具体的SQL语句，如：

---

[root@nineyou0-43 data]#mysqlbinlog--start-position=203 test.000004

/*!40019 SET@@session.max_insert_delayed_threads=0*/;

….

#090927 15:43:11 server id 1 end_log_pos 376 Query thread_id=188 exec_time=1 error_code=0

SET TIMESTAMP=1254037391/*!*/;

update t2 set username=upper(username)where id=1

/*!*/;

#at 376

#090927 15:43:11 server id 1 end_log_pos 403 Xid=1009

COMMIT/*!*/;

DELIMITER;

#End of log file

ROLLBACK/*added by mysqlbinlog*/;

/*!50003 SET COMPLETION_TYPE=@OLD_COMPLETION_TYPE*/;

---

## 对于SQL语句UPDATE t2 SET username=UPPER（username） WHERE id＝1，其记录的不是运行的SQL语句，而是更改后的结果。读者可能对这个表示不满，因为mysqlbinlog和Oracle LogMiner不同，解析出来的不是ROW格式下的SQL语句，因此mysqlbinlog的结果被"人为"地unreadable（难读）。

---

[root@nineyou0-43 data]#mysqlbinlog--start-position=1065 test.000004

/*!40019 SET@@session.max_insert_delayed_threads=0*/;

……

#at 1135

#at 1198

#090927 15:53:52 server id 1 end_log_pos 1198 Table_map:'member'.'t2'mapped to number 58

#090927 15:53:52 server id 1 end_log_pos 1378 Update_rows:table id 58 flags:STMT_END_F

BINLOG'

EBq/ShMBAAAAPwAAAK4EAAAAADoAAAAAAAAABm1lbWJlcgACdDIACgMPDw/+CgsPAQwKJAAoAEAA

/gJAAAAA

EBq/ShgBAAAAtAAAAGIFAAAQADoAAAAAAEACv////8A/AEAAAALYWxleDk5Dh5b3UEOXlvdSA3

Y2JiMzI1MmJhNmI3ZTljNDIyZmjNTMzNGGQyMjA1NAFNLacPAAAAAABjEnpxPBIAAAD8AQAAAtB

TEVYYTk4OFPVQQ5eW91IDdjYmIzMjUyYmE2YjdjdlOWM0MjJmYWM1MzM0ZDIyMDU0AU0tpw8AAAAA

AGMSenE8EgAA

'/*!*/;

#at 1378

#090927 15:53:52 server id 1 end_log_pos 1405 Xid=1110

COMMIT/*!*/;

DELIMITER;

#End of log file

ROLLBACK/*added by mysqlbinlog*/;

/*!50003 SET COMPLETION_TYPE=@OLD_COMPLETION_TYPE*/;

---

# 无论是哪种类型的SQL语句，要得到详细的信息，我们必须加上参数-v或-vv，两者的区别在于，前者比后者少-vv那部分-v没有的信息，其实通过对比就是-vv中的注释部分。

---

[root@nineyou0-43 data]#mysqlbinlog-vv--start-position=1065 test.000004

……

BINLOG'

EBq/ShMBAAAAPwAAAK4EAAAAADoAAAAAAAABm1lbWJlcgACdDIACgMPDw/+CgsPAQwKJAAoAEAA

/gJAAAAA

EBq/ShgBAAAAtAAAAGIFAAAQADoAAAAAAAEACv////8A/AEAAAALYWxleDk5Dh5b3UEOXlvdSSA3

Y2JiMzI1MmJhNmI3ZTljNDIyZmFjNTMzNGQyMjA1NAFNFNLacPAAAAAABjEnpxPBIAAAD8AQAAAtB

TEVYYOTk4OFlPVQQ5eW91IDdjYmYmIzMjUyYmE2YjdlOWM0MjJmYWM1MzM0ZDIyMDU0AU0tpw8AAAAA

AGMSenE8EgAA

'/*!*/;

###UPDATE member.t2

###WHERE

###@1=1/*INT meta=0 nullable=0 is_null=0*/

###@2='david'/*VARSTRING(36)meta=36 nullable=0 is_null=0*/

###@3='family'/*VARSTRING(40)meta=40 nullable=0 is_null=0*/

###@4='7cbb3252ba6b7e9c422fac5334d22054'/*VARSTRING(64)meta=64 nullable=0 is_null=0*/

###@5='M'/*STRING(2)meta=65026 nullable=0 is_null=0*/

###@6='2009:09:13'/*DATE meta=0 nullable=0 is_null=0*/

###@7='00:00:00'/*TIME meta=0 nullable=0 is_null=0*/

###@8=''/*VARSTRING(64)meta=64 nullable=0 is_null=0*/

###@9=0/*TINYINT meta=0 nullable=0 is_null=0*/

###@10=2009-08-11 16:32:35/*DATETIME meta=0 nullable=0 is_null=0*/

###SET

###@1=1/*INT meta=0 nullable=0 is_null=0*/

###@2='DAVID'/*VARSTRING(36)meta=36 nullable=0 is_null=0*/

###@3=family/*VARSTRING(40)meta=40 nullable=0 is_null=0*/

###@4='7cbb3252ba6b7e9c422fac5334d22054'/*VARSTRING(64)meta=64 nullable=0 is_null=0*/

###@5='M'/*STRING(2)meta=65026 nullable=0 is_null=0*/

###@6='2009:09:13'/*DATE meta=0 nullable=0 is_null=0*/

###@7='00:00:00'/*TIME meta=0 nullable=0 is_null=0*/

###@8=''/*VARSTRING(64)meta=64 nullable=0 is_null=0*/

###@9=0/*TINYINT meta=0 nullable=0 is_null=0*/

###@10=2009-08-11 16:32:35/*DATETIME meta=0 nullable=0 is_null=0*/

#at 1378

#090927 15:53:52 server id 1 end_log_pos 1405 Xid=1110

COMMIT/*!*/;

DELIMITER;

#End of log file

ROLLBACK/*added by mysqlbinlog*/;

/*!50003 SET COMPLETION_TYPE=@OLD_COMPLETION_TYPE*/;

---

从上mysqlbinlog解析出来的日志，对于一个简单的语句，比如update t2 set username=upper(username)where id=1这样的语句，也会记录所有字段原来的值，如果一个表有10W行数据，那么ROW格式的日志很可能会超过13MB。

## 3.3　查看套接字

如果你使用的是UNIX套接字方式访问MySQL，那么UNIX套接字文件所在的位置由参数□socket控制。一般来说，该套接字socket文件位于目录/tmp下，文件名为mysql.sock。

---

```
mysql□SHOW VARIABLES LIKE'socket'\G;

***************************1.row***************************

Variable_name:socket

Value:/tmp/mysql.sock

1 row in set(0.00 sec)
```

---

# 3.4　pid□□

□MySQL□□□□□□□□□□□□□□ID□□□□□□□□——□□□□□□pid□□□□□□□□□□□□□
pid_file□□□□□□□□□□□□□□□□□□□□□□□□□□□.pid□

---

mysql□show variables like'pid_file'\G;

***************************1.row***************************

Variable_name:pid_file

Value:/usr/local/mysql/data/xen-server.pid

1 row in set(0.00 sec)

---

# 3.5　视图的存储结构

由于MySQL的视图属于虚拟表，查看MySQL数据库系统的存储目录时，只能看到表示视图定义的frm文件，并没有单独的数据文件来存放视图的数据。

frm文件用来保存视图的定义信息，如果视图名为v_a，那么其对应的视图文件为v_a.frm。此视图文件是一个文本文件，可以直接使用cat命令查看其内容。

---

```
[root@xen-server test]#cat v_a.frm

TYPE=VIEW

query=select'test'.'a'.'b'AS'b'from'test'.'a'

md5=4eda70387716a4d6c96f3042dd68b742

updatable=1

algorithm=0

definer_user=root

definer_host=localhost

suid=2

with_check_option=0

timestamp=2010-08-04 07:23:36

create-version=1

source=select*from a

client_cs_name=utf8

connection_cl_name=utf8_general_ci

view_body_utf8=select'test'.'a'.'b'AS'b'from'test'.'a'
```

---

# 3.6 InnoDB体系结构

本节将从底层具体分析MySQL数据库的体系结构。在前面的几节中已经或多或少地提到了一些概念，本节将从InnoDB存储引擎的角度重新审视这些概念，同时还将介绍一些新的概念。

## 3.6.1 表空间

InnoDB存储引擎有一个被称为表空间（tablespace）的概念，它由一组数据文件组成，初始化的文件大小为10MB，名为ibdata1。我们将这个数据文件称为表空间文件（tablespace file），它由参数innodb_data_file_path定义，其基本的格式如下：

```
innodb_data_file_path=datafle_spec1[;datafle_spec2]...
```

我们可以通过多个文件组成一个表空间，同时制定文件的属性，如：

```
[mysqld]

innodb_data_file_path=/db/ibdata1:2000M;/dr2/db/ibdata2:2000M:autoextend
```

这里用/db/ibdata1和/dr2/db/ibdata2两个文件来组成表空间。注意，若这两个文件位于不同的磁盘上，磁盘的负载可能被平均，因此可以提高数据库的整体性能。两个文件的文件名之后都跟了属性，表示idbdata1的大小为2000MB，而ibdata2的大小为2000MB，如果用完了这2000MB，该文件可以自动增长（autoextend）。

设置innodb_data_file_path参数之后，所有基于InnoDB存储引擎的表的数据都会记录到该共享表空间中。若设置了参数innodb_file_per_table，则用户可以将每个基于InnoDB存储引擎的表产生一个独立表空间。独立表空间的命名规则为：表名.ibd。通过这样的方式，用户不用将所有数据都存放于默认的表空间中。下面这台MySQL数据库服务器设置了innodb_file_per_table，所以我们可以分别查看：

```
mysql＞SHOW VARIABLES LIKE'innodb_file_per_table'\G;

***************************1.row***************************

Variable_name:innodb_file_per_table
```

```
Value:ON

1 row in set(0.00 sec)

mysql〉system ls-lh/usr/local/mysql/data/member/*

-rw-r-----1 mysql mysql 8.7K 2009-02-24/usr/local/mysql/data/member/Profile.frm

-rw-r-----1 mysql mysql 1.7G 9〉25 11:13/usr/local/mysql/data/member/Profile.ibd

-rw-rw----1 mysql mysql 8.7K 9〉27 13:38/usr/local/mysql/data/member/t1.frm

-rw-rw----1 mysql mysql 17M 9〉27 13:40/usr/local/mysql/data/member/t1.ibd

-rw-rw----1 mysql mysql 8.7K 9〉27 15:42/usr/local/mysql/data/member/t2.frm

-rw-rw----1 mysql mysql 17M 9〉27 15:54/usr/local/mysql/data/member/t2.ibd
```

表Profile、t1、t2都采用InnoDB存储引擎，且由于设置了
innodb_file_per_table=ON，所以每张表都有.ibd文件与之对应，存放表中的数据。
那么这些数据是以什么样的方式组织的呢？下面通过BITMAP工具分析表中数据的组织方式，以帮助读者了
解3-1节介绍的InnoDB逻辑存储结构的关系。

图 3-1 InnoDB的表空间结构图

## 3.6.2 □□□□□□

□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□ib_logfile0□ib_logfile1□□□
□□□□MySQL□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□redo log file□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□
□□□□□□□□□□InnoDB□□□□□□□□□□□□□□

□□□□□□□□□□media failure□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□InnoDB□□□□□□□□□1□□□□□□□□□□□group□□□□□□□□□□□□□2□□□□□□□□□□□
□□□□ib_logfile0□ib_logfile1□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□mirrored log groups□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□1□□□□
□□□□□□□□□□□□□□□□□□□□□□2□□□□□□□□□□□2□□□□□□□□□□□□□□□□□□□□1□□□□
3-2□□□□□□□□□3□□□□□□□□□□□□□□□□□□□

图 3-2 日志文件组

日志组的重要参数包括如下几个：

☐innodb_log_file_size

☐innodb_log_files_in_group

☐innodb_mirrored_log_groups

☐innodb_log_group_home_dir

参数innodb_log_file_size用来设置重做日志文件的大小。在InnoDB1.2.x版本之前，重做日志文件总的大小不得大于等于4GB，而1.2.x版本将该限制扩大为了512GB。

参数innodb_log_files_in_group用来指定日志文件组中重做日志文件的数量，默认为2。参数innodb_mirrored_log_groups用来指定了日志镜像文件组的数量，默认为1，表示只有一个日志文件组，没有镜像。若磁盘本身已经做了高可用的方案，如磁盘阵列，那么可以不开启重做日志镜像的功能。最后，参数innodb_log_group_home_dir指定了日志文件组所在路径，默认为./，表示在MySQL数据库的数据目录下。以下显示了这些参数的当前值：

---

```
mysql∏SHOW VARIABLES LIKE'innodb%log%'\G;……

***************************4.row***************************

Variable_name:innodb_log_file_size

Value:5242880

***************************5.row***************************

Variable_name:innodb_log_files_in_group

Value:2

***************************6.row***************************

Variable_name:innodb_log_group_home_dir

Value:./

***************************7.row***************************

Variable_name:innodb_mirrored_log_groups

Value:1

7 rows in set(0.00 sec)
```

---

重做日志文件的大小设置对于InnoDB存储引擎的性能有着非常大的影响。一方面重做日志文件不能设置得太大，如果设置得很大，在恢复时可能需要很长的时间；另一方面又不能设置得太小了，否则可能导致一个事务的日志需要多次切换重做日志文件。此外，还会导致频繁地发生async checkpoint，导致性能的抖动。例如，用户可能会在错误日志中看到如下警告信息：

---

```
090924 11:39:44 InnoDB:ERROR:the age of the last checkpoint is 9433712,

InnoDB:which exceeds the log group capacity 9433498.

InnoDB:If you are using big BLOB or TEXT rows,you must set the
```

```
InnoDB:combined size of log files at least 10 times bigger than the

InnoDB:largest such row.

090924 11:40:00 InnoDB:ERROR:the age of the last checkpoint is 9433823,

InnoDB:which exceeds the log group capacity 9433498.

InnoDB:If you are using big BLOB or TEXT rows,you must set the

InnoDB:combined size of log files at least 10 times bigger than the

InnoDB:largest such row.

090924 11:40:16 InnoDB:ERROR:the age of the last checkpoint is 9433645,

InnoDB:which exceeds the log group capacity 9433498.

InnoDB:If you are using big BLOB or TEXT rows,you must set the

InnoDB:combined size of log files at least 10 times bigger than the

InnoDB:largest such row.
```

---

□□□□□□□InnoDB:ERROR:the age of the last checkpoint is 9433645□InnoDB:which exceeds the log group capacity 9433498□□□□□□□□□□□□□capacity□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□innodb buffer pool□□□□□□□□flush list□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□InnoDB□MyISAM□Heap□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□STATEMENT□□ROW□□□□□□MIXED□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□Page□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□redo entry□□□□□□□□□□□□□□□□□

□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB 1.2.x□□□□□□□□□□□□51□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□3-2□□□□□□□□□□□□□□

表 3-2 重做日志条目结构

| redo_log_type | space | page_no | redo_log_body |
|---|---|---|---|

表格3-2中各个部分的含义如下4个部分组成：

☐redo_log_type：占用1字节，表示重做日志的类型。

☐space：表示表空间的ID，但是在存储空间中占用的空间可能小于4字节。

☐page_no：表示页号，同样可能被压缩存储。

☐redo_log_body：表示每个重做日志的数据部分，恢复时会调用相应的函数进行解析。

前面2小节介绍了一条重做日志的基本格式，但是一条记录往往会同时修改redo log buffer。其修改方式同一个事务对应多条重做日志，如图3-3所示。每个数据库对应的重做日志

图 3-3 重做日志文件组

重做日志文件以块为单位，每块大小512字节进行保存。同样，在每次写入重做日志文件时，都需要调用一次文件系统的sync操作，因为重做日志缓冲先写入文件系统缓存。为了确保重做日志写入磁盘，必须进行一次doublewrite。

在之前的章节中，我们已经知道从日志缓冲写入磁盘操作，在主线程（master thread）中每秒会进行一次。下面分析从重做日志缓冲往磁盘写入重做日志文件的操作。此外，参数innodb_flush_log_at_trx_commit用来控制重做日志刷新到磁盘的策略。

参数innodb_flush_log_at_trx_commit的有效值有0、1、2。0代表当提交事务时，并不将事务的重做日志写入磁盘上的日志文件，而是等待主线程每秒的刷新。1和2不同的是：1表示在执行commit时将重做日志缓冲同步写到磁盘，即伴有fsync的调用。2表示将重做日志异步写到磁盘，即写到文件系统的缓存中。因此不能完全保证在执行commit时肯定会写入重做日志文件，只是有这个动作发生。

□□□□□□□□□□ACID□□□□□□□□□□innodb_flush_log_at_trx_commit□□□□1□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□0□2□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□2□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 3.7　□□

□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□
□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□
□□□□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□point in time□□□□□□□□□replication□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL 5.1□□□□□□□□□□□□□
STATEMENT□ROW□MIX□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□

# 第4章 锁

本章将对 InnoDB 存储引擎中锁的实现和其内部机制做深入的分析。希望通过对本章知识的学习，读者能将各种锁机制应用于实际的生产环境，避免不必要的错误，更好地使用数据库。

# 4.1　主键的选择

由InnoDB存储引擎的逻辑存储结构可知，数据总是按照主键顺序进行存放，即index organized table，在InnoDB存储引擎中如果没有显式地创建Primary Key，那么存储引擎会按照如下方式选择或创建主键。首先InnoDB存储引擎会判断表中是否有非空的唯一索引（Unique NOT NULL），如果有，则该列即为主键。

如果不符合上述条件，InnoDB存储引擎会自动创建一个6字节大小的指针。

当表中有多个非空唯一索引时，InnoDB存储引擎会选择表中第一个定义的非空唯一索引为主键。这里需要非常注意的是，主键的选择是根据定义索引的顺序，而不是建表时列的顺序。下面我们通过例子来分析这个问题。

```
mysql>CREATE TABLE z(

->a INT NOT NULL,

->b INT NULL,

->c INT NOT NULL,

->d INT NOT NULL,

->UNIQUE KEY(b),

->UNIQUE KEY(d),UNIQUE KEY(c));

Query OK,0 rows affected(0.02 sec)

mysql>INSERT INTO z SELECT 1,2,3,4;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql>INSERT INTO z SELECT 5,6,7,8;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql>INSERT INTO z SELECT 9,10,11,12;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0
```

在我们创建的表z中，a、b、c、d列内容和b、c、d列内容完全一样，并且都不是NULL。由此我们可以知道主键是怎么生成的了。接着继续看下面的SQL语句，看看运行后的结果。

---

mysql＞SELECT a,b,c,d,_rowid FROM z;

+---+------+----+----+--------+

|a|b|c|d|_rowid|

+---+------+----+----+--------+

|1|2|3|4|4|

|5|6|7|8|8|

|9|10|11|12|12|

+---+------+----+----+--------+

3 rows in set(0.00 sec)

---

# _rowid可以显示表的主键，因此通过上述代码可以发现，表z的主键为列d的值。这个例子再次强调：_rowid只能用于查看单个列为主键的情况，对于多列组成的主键就显得无能为力了。列d是被自动创建的，因此是InnoDB存储引擎的主键。

# 下面的例子显示了_rowid不能用于多列主键的InnoDB存储引擎表的主键查询。

---

mysql＞CREATE TABLE a(

-＞a INT,

-＞b INT,

-＞PRIMARY KEY(a,b)

-＞)ENGINE=InnoDB;

Query OK,0 rows affected(0.03 sec)

mysql＞INSERT INTO a SELECT 1,1;

Query OK,1 row affected(0.01 sec)

Records:1 Duplicates:0 Warnings:0

mysql＞SELECT a,_rowid FROM a;

ERROR 1054(42S22):Unknown column'_rowid'in'field list'

## 4.2　InnoDB逻辑存储结构

从InnoDB的逻辑存储结构看，所有数据都被逻辑地存放在一个空间中，称为表空间（tablespace）。表空间又由段（segment）、区（extent）、页（page）组成。页在一些文档中有时也称为块（block），InnoDB存储引擎的逻辑存储结构大致如图4-1所示。

## Tablespace

Leaf node segment

Non-Leaf node segment

Rollback segment

## Segment

| Extent | Extent |
|--------|--------|
| Extent | Extent |

## Extent

## Page

| Row | Row |
|-----|-----|
| Row | Row |
| Row | |

## Row

Trx id

Roll Pointer

Roll Pointer

| Col 1 | ... | Col n |

图 4-1 InnoDB表的存储方式

# 4.2.1 表空间

表空间可以看做是InnoDB存储引擎逻辑结构的最高层，所有的数据都存放在表空间中。然而在默认情况下InnoDB存储引擎有一个共享表空间ibdata1，即所有数据都存放在这个表空间内。如果用户启用了参数innodb_file_per_table，则每张表内的数据可以单独放到一个表空间内。

如果启用了innodb_file_per_table的参数，需要注意每张表的表空间内存放的只是数据、索引和插入缓冲Bitmap页，其他类的数据，如回滚（undo）信息，插入缓冲索引页、系统事务信息，二次写缓冲（Double write buffer）等还是存放在原来的共享表空间内。这同时也说明了另一个问题：即使在启用了参数innodb_file_per_table之后，共享表空间还是会不断地增加其大小。下面用一个例子来显示在启用innodb_file_per_table后的innodb存储引擎存储分布情况，首先对当前数据库下的表mytest做一次完整的操作，然后看表空间的存储分布：

---

```
mysql＞SHOW VARIABLES LIKE'innodb_file_per_table'\G;

***************************1.row***************************

Variable_name:innodb_file_per_table

Value:ON

1 row in set(0.00 sec)

mysql＞system ls-lh/usr/local/mysql/data/ibdata*

-rw-rw----1 mysql mysql 58M Mar 11 13:58/usr/local/mysql/data/ibdata1
```

---

可以看到，共享表空间ibdata1的大小为58MB。这里由于共享的undo信息，插入缓冲等1信息的表空间mytest还是存放在默认的共享InnoDB系统表空间中。

---

```
mysql＞SET autocommit=0;

Query OK,0 rows affected(0.00 sec)

mysql＞UPDATE mytest SET salary=0;

Query OK,2844047 rows affected(19.47 sec)

Rows matched:2844047 Changed:2844047 Warnings:0

mysql＞system ls-lh/usr/local/mysql/data/ibdata*

-rw-rw----1 mysql mysql 114M Mar 11 14:00/usr/local/mysql/data/ibdata1
```

将所有记录的列值更新为0，但是仍然看不到实际操作过程中磁盘对undo log的操作。在事务提交后，通过commit、rollback去查看，还是会看到对undo的操作。update mytest set salary=0之后再查看，通过命令发现ibdata1文件变为了114MB，这个增加的部分即为操作过程中所产生的undo信息。

接着我们进行回滚，k时刻后如果rollback，ibdata1文件也不会因此而发生磁盘容量的收缩58MB，这就证明了前面所说的，事务回滚通过undo进行。

mysql>ROLLBACK;

Query OK,0 rows affected(0.00 sec)

mysql>system ls-lh/usr/local/mysql/data/ibdata*

-rw-rw----1 mysql mysql 114M Mar 11 14:00/usr/local/mysql/data/ibdata1

从"缩小"的例子中可以看出，虽然114MB的InnoDB存储引擎不会因为rollback而去收缩这个文件。虽然InnoDB不会回收这些空间，但是会自动判断这些undo信息是否还需要，并将这些空间标记为可用空间，供下次undo使用。

master thread每10秒会执行一次full purge操作，因此很有可能的情况是，用户再次执行上述的UPDATE语句后，会发现ibdata1文件不会再次增大了，这就是这个道理。

此外，用户可以通过python脚本py_innodb_page_info工具来统计表空间中各页的类型和数量，该工具可以在code.google.com上搜索david-mysql-tools进行下载，运行结果如下。

[root@nineyou0-43 py]#python py_innodb_page_info.py/usr/local/mysql/data/ibdata1

Total number of page:83584:

Insert Buffer Free List:204

Freshly Allocated Page:5467

Undo Log Page:38675

File Segment inode:4

B-tree Node:39233

File Space Header:1

□□□□□□83 584□□□□□□□□□□□□□□□□□□204□□□5467□□□□□38 675□undo□□39 233□□□□□□□□□□□□□□□□-v□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 4.2.2 □

图4-1显示了一张表中数据的组织方式，其中的数据都存放在一个聚集索引中。因为
InnoDB存储引擎表是索引组织表（index organized），即数据即索引，索引即数据。而
聚集索引的B+树叶子节点就是图4-1中Leaf node segment，在聚集索引的B+树非叶子节点
就是图4-1中Non-leaf node segment。此外，还有一类特殊的段用来管理大数据类型，即
□。

对InnoDB存储引擎表来说，一个索引会产生两个段。而对于每个段的管理，DBA其实并不需要太过关心。相对于Oracle数据库对段的管理来说，由于ASSM的出现，现在也不需要DBA过多的干预。

# 4.2.3 □

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□1MB□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□4□5□□□□□□□□□□□InnoDB□□□□□□□□□□□16KB□□□□□□□□□□□64□□□□□□□□

InnoDB 1.0.x□□□□□□□□□□□□□□□□□□□□□□□□□□□KEY_BLOCK_SIZE□□□□2K□4K□8K□□□□□□□□□□□□□□□□□□□□512□256□128□

InnoDB 1.2.x□□□□□□□□innodb_page_size□□□□□□□□□□□□□□□□□□□□□□□4K□8K□□□□□□□□□□□□□□□□□□□□□□□□□256□128□□□□□□□□□□□□□□□□□□□□□□□□□1M□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□innodb_file_per_talbe□□□□□□□□□□□□□□□96KB□□□□□64□□□□□□□□□□□□□□□□□□□1MB□□□□□□□□□□□□□□□□□□□□□□32□□□□□□□□□□fragment page□□□□□□□□□□□□□□□□□□□□□□64□□□□□□□□□□□□□□□□□□□□□□□□□□□□□undo□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□

---

```
mysql□CREATE TABLE t1(

-□col1 INT NOT NULL AUTO_INCREMENT,

-□col2 VARCHAR(7000),

-□PRIMARY KEY(col1))ENGINE=InnoDB;

mysql□system ls-lh/usr/local/mysql/data/test/t1.ibd;

-rw-rw----1 mysql mysql 96K 10□12 14:59/usr/local/mysql/data/test/t1.ibd
```

---

□□□□SQL□□□□□□□t1□□□□col2□□□□□VARCHAR□7000□□□□□□□□□□□□□□□□□□□□□2□□□□□□□□ls□□□□□□□□□□□□□□□□□t1□□□□□□□□□□□□□96KB□□□□□□□□□SQL□□□□

---

```
mysql□INSERT t1 SELECT NULL,REPEAT('a',7000);

Query OK,1 row affected(0.04 sec)
```

Records:1 Duplicates:0 Warnings:0

mysql□INSERT into t1 SELECT NULL,REPEAT('a',7000);

Query OK,1 row affected(0.01 sec)

Records:1 Duplicates:0 Warnings:0

mysql□system ls-lh/usr/local/mysql/data/test/t1.ibd;

-rw-rw----1 mysql mysql 96K 10□12 16:24/usr/local/mysql/data/test/t1.ibd

## 这里又插入了一行记录。我们再次通过工具来观察表空间文件，并用上一节提到的 py_innodb_page_info工具来进行分析，可得：

[root@nineyou0-43 py]#./py_innodb_page_info.py-v/usr/local/mysql/data/test/t1.ibd

page offset 00000000,page type□File Space Header□

page offset 00000001,page type□Insert Buffer Bitmap□

page offset 00000002,page type□File Segment inode□

page offset 00000003,page type□B-tree Node□,page level□0000□

page offset 00000000,page type□Freshly Allocated Page□

page offset 00000000,page type□Freshly Allocated Page□

Total number of page:6:

Freshly Allocated Page:2

Insert Buffer Bitmap:1

File Space Header:1

B-tree Node:1

File Segment inode:1

## 通过-v详细信息我们可以观察到，目前数据页（page offset为3）的这个页的页类型为（page level）的值仍然为0，代表数据页仍然只有一个页。然而我们插入的数据的确变大了，因此这是怎么回事呢？我们继续插入数据。

mysql□INSERT into t1 SELECT NULL,REPEAT('a',7000);

Query OK,1 row affected(0.01 sec)

Records:1 Duplicates:0 Warnings:0

[root@nineyou0-43 py]#./py_innodb_page_info.py-v/usr/local/mysql/data/test/t1.ibd

page offset 00000000,page type□File Space Header□

page offset 00000001,page type□Insert Buffer Bitmap□

page offset 00000002,page type□File Segment inode□

page offset 00000003,page type□B-tree Node□,page level□0001□

page offset 00000004,page type□B-tree Node□,page level□0000□

page offset 00000005,page type□B-tree Node□,page level□0000□

Total number of page:6:

Insert Buffer Bitmap:1

File Space Header:1

B-tree Node:3

File Segment inode:1

---

# 从结果中可以看到page offset为3的页的page level已经是1，不再是0了。这时B+树的高度已经变为2了，B-tree Node□

## 当表的记录为，，，，□其实就是60万条记录时（可以看到表t1的记录为63万条，在32万条前，此时的高度已经为2了，故确保此时高度为2的条件可以满足

---

mysql□DELIMITER//

mysql□CREATE PROCEDURE load_t1(count INT UNSIGNED)

-□BEGIN

-□DECLARE s INT UNSIGNED DEFAULT 1;

-□DECLARE c VARCHAR(7000)DEFAULT REPEAT('a',7000);

-□WHILE s□=count DO

-□INSERT INTO t1 SELECT NULL,c;

-□SET s=s+1;

-□END WHILE;

-□END;

-□//

Query OK,0 rows affected(0.04 sec)

mysql□DELIMITER;

mysql□CALL load_t1(60);

Query OK,1 row affected(1.59 sec)

mysql□SELECT COUNT(*)FROM t1\G;

***************************1.row***************************

count(*):63

1 row in set(0.00 sec)1 row in set(0.00 sec)

mysql□system ls-lh/usr/local/mysql/data/test/t1.ibd;

-rw-rw----1 mysql mysql 576K 10□12 16:56/usr/local/mysql/data/test/t1.ibd

---

# 从上面的例子可以看到，63个页还是小于1MB的，因此还是在碎片区中进行空间的分配。由于63个页是从第0个页开始记录的，用py_innodb_page_info工具来分析表空间（t.ibd），可得：

---

[root@nineyou0-43 py]#./py_innodb_page_info.py-v/usr/local/mysql/data/test/t1.ibd

page offset 00000000,page type□File Space Header□

page offset 00000001,page type□Insert Buffer Bitmap□

page offset 00000002,page type□File Segment inode□

page offset 00000003,page type□B-tree Node□,page level□0001□

page offset 00000004,page type□B-tree Node□,page level□0000□

page offset 00000005,page type□B-tree Node□,page level□0000□

page offset 00000006,page type□B-tree Node□,page level□0000□

page offset 00000007,page type□B-tree Node□,page level□0000□

page offset 00000008,page type□B-tree Node□,page level□0000□

page offset 00000009,page type□B-tree Node□,page level□0000□

page offset 0000000a,page type□B-tree Node□,page level□0000□

page offset 0000000b,page type□B-tree Node□,page level□0000□

page offset 0000000c,page type□B-tree Node□,page level□0000□

page offset 0000000d,page type□B-tree Node□,page level□0000□

page offset 0000000e,page type【B-tree Node】,page level【0000】

page offset 0000000f,page type【B-tree Node】,page level【0000】

page offset 00000010,page type【B-tree Node】,page level【0000】

page offset 00000011,page type【B-tree Node】,page level【0000】

page offset 00000012,page type【B-tree Node】,page level【0000】

page offset 00000013,page type【B-tree Node】,page level【0000】

page offset 00000014,page type【B-tree Node】,page level【0000】

page offset 00000015,page type【B-tree Node】,page level【0000】

page offset 00000016,page type【B-tree Node】,page level【0000】

page offset 00000017,page type【B-tree Node】,page level【0000】

page offset 00000018,page type【B-tree Node】,page level【0000】

page offset 00000019,page type【B-tree Node】,page level【0000】

page offset 0000001a,page type【B-tree Node】,page level【0000】

page offset 0000001b,page type【B-tree Node】,page level【0000】

page offset 0000001c,page type【B-tree Node】,page level【0000】

page offset 0000001d,page type【B-tree Node】,page level【0000】

page offset 0000001e,page type【B-tree Node】,page level【0000】

page offset 0000001f,page type【B-tree Node】,page level【0000】

page offset 00000020,page type【B-tree Node】,page level【0000】

page offset 00000021,page type【B-tree Node】,page level【0000】

page offset 00000022,page type【B-tree Node】,page level【0000】

page offset 00000023,page type【B-tree Node】,page level【0000】

Total number of page:36:

Insert Buffer Bitmap:1

File Space Header:1

B-tree Node:33

File Segment inode:1

我们发现，B-tree Node一共有33个，除去一个page level为1的非叶子节点，一共有32个page level为0的叶子节点，也就是说，一共有32个叶子节点用来存放我们插入的数据。由于我们

# 即插入64条记录，正好是一个簇。我们通过脚本工具观察插入这些数据后的表空间分布情况。

---

mysql＞CALL load_t1(1);

Query OK,1 row affected(0.10 sec)

mysql＞system ls-lh/usr/local/mysql/data/test/t1.ibd;

-rw-rw----1 mysql mysql 2.0M 10月12 17:02/usr/local/mysql/data/test/t1.ibd

---

# 其结果显示共有32个碎片页被使用（我们同样可以用前面小节讲过的py_innodb_page_info工具来进行分析，t1.ibd的后续碎片页属于类型Freshly Allocated Page的页）

---

[root@nineyou0-43 test2]#～/py/py_innodb_page_info.py t1.ibd-v

page offset 00000000,page type＜File Space Header＞

page offset 00000001,page type＜Insert Buffer Bitmap＞

page offset 00000002,page type＜File Segment inode＞

page offset 00000003,page type＜B-tree Node＞,page level＜0001＞

page offset 00000004,page type＜B-tree Node＞,page level＜0000＞

page offset 00000005,page type＜B-tree Node＞,page level＜0000＞

page offset 00000006,page type＜B-tree Node＞,page level＜0000＞

page offset 00000007,page type＜B-tree Node＞,page level＜0000＞

page offset 00000008,page type＜B-tree Node＞,page level＜0000＞

page offset 00000009,page type＜B-tree Node＞,page level＜0000＞

page offset 0000000a,page type＜B-tree Node＞,page level＜0000＞

page offset 0000000b,page type＜B-tree Node＞,page level＜0000＞

page offset 0000000c,page type＜B-tree Node＞,page level＜0000＞

page offset 0000000d,page type＜B-tree Node＞,page level＜0000＞

page offset 0000000e,page type＜B-tree Node＞,page level＜0000＞

page offset 0000000f,page type＜B-tree Node＞,page level＜0000＞

page offset 00000010,page type＜B-tree Node＞,page level＜0000＞

page offset 00000011,page type＜B-tree Node＞,page level＜0000＞

page offset 00000012,page type「B-tree Node」,page level「0000」

page offset 00000013,page type「B-tree Node」,page level「0000」

page offset 00000014,page type「B-tree Node」,page level「0000」

page offset 00000015,page type「B-tree Node」,page level「0000」

page offset 00000016,page type「B-tree Node」,page level「0000」

page offset 00000017,page type「B-tree Node」,page level「0000」

page offset 00000018,page type「B-tree Node」,page level「0000」

page offset 00000019,page type「B-tree Node」,page level「0000」

page offset 0000001a,page type「B-tree Node」,page level「0000」

page offset 0000001b,page type「B-tree Node」,page level「0000」

page offset 0000001c,page type「B-tree Node」,page level「0000」

page offset 0000001d,page type「B-tree Node」,page level「0000」

page offset 0000001e,page type「B-tree Node」,page level「0000」

page offset 0000001f,page type「B-tree Node」,page level「0000」

page offset 00000020,page type「B-tree Node」,page level「0000」

page offset 00000021,page type「B-tree Node」,page level「0000」

page offset 00000022,page type「B-tree Node」,page level「0000」

page offset 00000023,page type「B-tree Node」,page level「0000」

page offset 00000000,page type「Freshly Allocated Page」

……

page offset 00000000,page type「Freshly Allocated Page」

page offset 00000000,page type「Freshly Allocated Page」

page offset 00000000,page type「Freshly Allocated Page」

Total number of page:128:

Freshly Allocated Page:91

Insert Buffer Bitmap:1

File Space Header:1

B-tree Node:34

File Segment inode:1

## 4.2.4　□

□□□□□□□□□□InnoDB□□□□Page□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□
□□□□□□InnoDB□□□□□□□□□□□□□□□□16KB□□□□InnoDB 1.2.x□□□□□□□
□□□□□□innodb_page_size□□□□□□□□□□4K□8K□16K□□□□□□□□□□□□□□□□
□□□□□□innodb_page_size□□□□□□□□□□□□□□□□□□□□mysqldump□□□□□
□□□□□□□□□□□□□

□InnoDB□□□□□□□□□□□□□□□□□□

□□□□□□B-tree Node□

□undo□□undo Log Page□

□□□□□System Page□

□□□□□□□Transaction system Page□

□□□□□□□□□Insert Buffer Bitmap□

□□□□□□□□□□□□Insert Buffer Free List□

□□□□□□□□□□□□□□□Uncompressed BLOB Page□

□□□□□□□□□□□□□□compressed BLOB Page□

## 4.2.5 □

InnoDB存储引擎是基于行（row-oriented）的数据库，因此对于较大的列，如果将这部分数据存放在页的某部分，16KB/2-200，即超过了7992行的情况下。数据库基于行（row-oriented）的存储引擎，而不是基于列（column-oriented）的。MySQL infobright存储引擎是基于列的，因此在一些应用下，基于列的数据库可以更好地运行SQL查询。同样，一些数据库也是基于列的，如Sybase IQ、Google Big Table。基于列的数据库，适合于数据仓库下的分析型应用，而对于联机事务处理则并不适合。

# 4.3 InnoDB行记录格式

InnoDB存储引擎和大多数数据库一样（如Oracle、Microsoft SQL Server数据库等），记录是以行的形式存储的。这意味着页中保存着表中一行行的数据。在InnoDB 1.0.x版本之前，InnoDB存储引擎提供了Compact和Redundant两种格式来存放行记录数据，这也是目前使用最多的一种格式。Redundant格式是为兼容之前版本而保留的。如果阅读过InnoDB的源代码，用户会发现两种格式都被称为PHYSICAL RECORD（NEW STYLE）和PHYSICAL RECORD（OLD STYLE）。由于从MySQL 5.1版本开始，默认设置为Compact行格式。用户可以通过命令SHOW TABLE STATUS LIKE'table_name'来查看当前表使用的行格式，其中row_format属性表示当前所使用的行记录结构类型。

---

```
mysql：SHOW TABLE STATUS like'mytest%'\G;

***************************1.row***************************

Name:mytest

Engine:InnoDB

Version:10

Row_format:Compact

Rows:6

Avg_row_length:2730

Data_length:16384

Max_data_length:0

Index_length:0

Data_free:0

Auto_increment:NULL

Create_time:2009-03-17 13:33:50

Update_time:NULL

Check_time:NULL

Collation:latin1_swedish_ci

Checksum:NULL

Create_options:
```

Comment:

*************************2.row*************************

Name:mytest2

Engine:InnoDB

Version:10

Row_format:Redundant

Rows:0

Avg_row_length:0

Data_length:16384

Max_data_length:0

Index_length:0

Data_free:0

Auto_increment:NULL

Create_time:2009-03-17 13:57:23

Update_time:NULL

Check_time:NULL

Collation:latin1_swedish_ci

Checksum:NULL

Create_options:row_format=REDUNDANT

Comment:

2 rows in set(0.00 sec)

可以看到，表mytest是Compact的行格式，mytest2是Redundant的行格式。接着分析表空间文件，由于表的内容非常简单，通过前面的知识应该可以分析，这里不再具体分析表空间文件，有兴趣的读者可以用之前的py_innodb_page_info工具，并通过文本编辑软件来具体分析这两种不同行记录格式的区别。

# 4.3.1　Compact行记录格式

Compact行记录是在MySQL 5.0中引入的，其设计目标是高效地存储数据。简单来说，一个页中存放的行数据越多，其性能就越高。图4-2显示了Compact行记录的存储方式：

| 变长字段长度列表 | NULL标志位 | 记录头信息 | 列1数据 | 列2数据 | …… |
|---|---|---|---|---|---|

图 4-2 Compact行记录格式

由图4-2可以看到，Compact行记录格式的首部是一个非NULL变长字段长度列表，并且其是按照列的顺序逆序放置的，其长度为：

□若列的长度小于255字节，用1字节表示；

□若大于255字节，用2字节表示。

变长字段的长度最大不可以超过2字节，这是因为MySQL数据库中VARCHAR类型的最大长度限制为65535。第二个部分是NULL标志位，该位指示了该行数据中是否有NULL值，有则用1表示。该部分所占的字节应该为1字节。接下来的部分是记录头信息（record header），固定占用5字节（40位），每位的含义见表4-1。

表 4-1　Compact 记录头信息

| 名　　称 | 大小 (bit) | 描　　述 |
|---|---|---|
| () | 1 | 未知 |
| () | 1 | 未知 |
| deleted_flag | 1 | 该行是否已被删除 |
| min_rec_flag | 1 | 为 1，如果该记录是预先被定义为最小的记录 |
| n_owned | 4 | 该记录拥有的记录数 |
| heap_no | 13 | 索引堆中该条记录的排序记录 |
| record_type | 3 | 记录类型，000 表示普通，001 表示 B+ 树节点指针，010 表示 Infimum，011 表示 Supremum，1xx 表示保留 |
| next_record | 16 | 页中下一条记录的相对位置 |
| Total | 40 | |

需要注意的是，不同于每行的数据需要特殊处理，对于NULL值的列，列的长度也必须为NULL，而不是NULL的列，其长度取决于该列的数据类型，而主键、事务ID、回滚指针这三列，6字节、7字节的值是由InnoDB自动生成的，而主键这里是6字节的rowid值。

# □□□□□□□□□□□□□□□Compact□□□□□□□□□□□

```
mysql□CREATE TABLE mytest(

-□t1 VARCHAR(10),

-□t2 VARCHAR(10),

-□t3 CHAR(10),

-□t4 VARCHAR(10)

-□)ENGINE=INNODB CHARSET=LATIN1 ROW_FORMAT=COMPACT;

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO mytest

-□VALUES('a','bb','bb','ccc');

Query OK,1 row affected(0.01 sec)

mysql□INSERT INTO mytest

-□VALUES('d','ee','ee','fff');

Query OK,1 row affected(0.00 sec)

mysql□INSERT INTO mytest

-□VALUES('d',NULL,NULL,'fff');

Query OK,1 row affected(0.00 sec)

mysql□SELECT*FROM mytest\G;

***************************1.row**************************

t1:a

t2:bb

t3:bb

t4:ccc

***************************2.row**************************

t1:d

t2:ee

t3:ee

t4:fff

***************************3.row**************************
```

t1:d

t2:NULL

t3:NULL

t4:fff

3 rows in set(0.00 sec)

---

上面创建了一个名为mytest的表，共有4个列t1、t2、t4均为VARCHAR类型，列t3为固定长度CHAR类型。插入了3条记录。我们现在打开表空间文件mytest.ibd（启用了innodb_file_per_table，如果没有启用，应该打开共享表空间文件ibdata1）。

用Windows的读者可以使用工具软件UltraEdit打开该文件，用Linux的读者可以使用hexdump-C-v mytest.ibd＞mytest.txt命令，之后用文本编辑器打开mytest.txt文件。mytest.txt文件的内容看起来如下：

---

0000c070 73 75 70 72 65 6d 75 6d 03 02 01 00 00 00 10 00|supremum........|

0000c080 2c 00 00 00 2b 68 00 00 00 00 00 06 05 80 00 00|,...+h..........|

0000c090 00 32 01 10 61 62 62 62 62 20 20 20 20 20 20 20|.2..abbbb|

0000c0a0 20 63 63 63 03 02 01 00 00 00 18 00 2b 00 00 00|ccc........+...|

0000c0b0 2b 68 01 00 00 00 00 06 06 80 00 00 00 32 01 10|+h...........2..|

0000c0c0 64 65 65 65 65 20 20 20 20 20 20 20 20 66 66 66|deeeefff|

0000c0d0 03 01 06 00 00 20 ff 98 00 00 00 2b 68 02 00 00|.........+h...|

0000c0e0 00 00 06 07 80 00 00 00 32 01 10 64 66 66 66 00|.......2..dfff.|

---

# 我们来看0000c078开始的这一行，即第一行被插入的记录：

---

03 02 01/*变长字段长度列表，逆序*/

00/*NULL标志位，第一行没有NULL值*/

00 00 10 00 2c/*Record Header，固定5字节长度*/

00 00 00 2b 68 00/*RowID InnoDB自动创建，6字节*/

00 00 00 00 06 05/*TransactionID*/

80 00 0000 32 01 10/*Roll Pointer*/

61/*第1行数据'a'*/

62 62/*第2行数据'bb'*/

62 62 20 20 20 20 20 20 20 20/*第3行数据'bb'*/

63 63 63/*第4行数据'ccc'*/

---

我们把上面这段数据中表示列内容的数据给抽了出来，可以看到数据库存储的确实就是03 02 01，而不是01 02 03。接着数据就是InnoDB自己添加的TransactionID和Roll Pointer。我们也可以看到对于CHAR类型的列，它会在内容不足的地方添加0x20来进行填充。

另外我们刚才在Record Header的解析中，看到下一条数据的偏移next_recorder是0x2c，我们在当前数据的起始位置加上偏移量0x2c，也就是下面这段数据。这一部分也是InnoDB添加的三个隐藏列以及相关的数据内容，我们继续来进行解析。

我们在这段数据中看到RowID，也就是说这个表没有主键。由于有两个字段是可以为空的，所以它有NULL值列表来记录。

---

03 01/*可以为空的两个字段的长度*/

06/*NULL值列表，表示有NULL值*/

00 00 20 ff 98/*Record Header*/

00 00 00 2b 68 02/*RowID*/

00 00 00 00 06 07/*TransactionID*/

80 00 00 00 32 01 10/*Roll Pointer*/

64/*第1行数据'd'*/

66 66 66/*第4行数据'fff'*/

---

由于这行有NULL值，所以NULL值列表的值为00，而06换成二进制就是00000110，第1个字段和第2个字段第3个字段是可以为NULL的字段，所以它们对应二进制从低位往高位数分别为1，而第4个NULL值列表的位顺序是颠倒的。由于是CHAR类型而非VARCHAR类型，在compact模式下NULL值列表是不会记录它们的。

## 4.3.2　Redundant行记录格式

Redundant是MySQL 5.0版本前InnoDB的行记录存储方式，MySQL 5.0支持Redundant是为了兼容之前版本的页格式。Redundant行记录格式如图4-3所示。



| 字段长度偏移列表 | 记录头信息 | 列1数据 | 列2数据 | 列3数据 | …… |

图　4-3　Redundant行记录格式

从图4-3可以看到，不同于Compact行记录格式，Redundant行记录格式的第一个部分是每一列长度偏移列表，同样也是逆序存放的。若列的长度小于255字节，用1字节表示；若大于255字节，用2字节表示。第二个部分为记录头信息（record header），不同于Compact行记录格式，Redundant行记录格式的头信息固定占用6字节（48位），每位的含义见表4-2。从表4-2中可以看到，n_fields值代表一行中列的数量，占用10位。同时这也很好地解释了为什么MySQL数据库一行支持最多的列为1023。另一个需要注意的值为1byte_offs_flags，该值定义了偏移列表占用1字节还是2字节。而最后的部分就是实际存储的每个列的数据了。

## 表 4-2 Redundant 记录头信息

| 名　　称 | 大小 (bit) | 描　　述 |
|---|---|---|
| () | 1 | 未知 |
| () | 1 | 未知 |
| deleted_flag | 1 | 该行是否已被删除 |
| min_rec_flag | 1 | 为 1，如果该记录是预先被定义为最小的记录 |
| n_owned | 4 | 该记录拥有的记录数 |
| heap_no | 13 | 索引堆中该条记录的索引号 |
| n_fields | 10 | 记录中列的数量 |
| 1byte_offs_flag | 1 | 偏移列表为 1 字节还是 2 字节 |
| next_record | 16 | 页中下一条记录的相对位置 |
| Total | 48 | |

我们创建和第 4.3.1 节中表 mytest 一样的表，存储引擎为 Redundant，表 mytest2，

mysql　CREATE TABLE mytest2

-　ENGINE=InnodB ROW_FORMAT=Redundant

-□AS

-□SELECT*FROM mytest;

Query OK,3 rows affected(0.00 sec)

Records:3 Duplicates:0 Warnings:0

mysql□SHOW TABLE STATUS LIKE'mytest2'\G;

***************************1.row***************************

Name:mytest2

Engine:InnoDB

Version:10

Row_format:Redundant

Rows:3

Avg_row_length:5461

Data_length:16384

Max_data_length:0

Index_length:0

Data_free:0

Auto_increment:NULL

Create_time:2009-03-18 15:49:42

Update_time:NULL

Check_time:NULL

Collation:latin1_swedish_ci

Checksum:NULL

Create_options:row_format=REDUNDANT

Comment:

1 row in set(0.00 sec)

mysql□SELECT*FROM mytest2\G;

***************************1.row***************************

t1:a

t2:bb

t3:bb

t4:ccc

***************************2.row***************************

t1:d

t2:ee

t3:ee

t4:fff

***************************3.row***************************

t1:d

t2:NULL

t3:NULL

t4:fff

3 rows in set(0.00 sec)

---

# 我们仍然使用row_format＝＝Redundant重新来使用hexdump方法分析mytest2.ibd文件，并把它存入mytest2.txt里，重要的几行内容如下：

---

```
0000c070 08 03 00 00 73 75 70 72 65 6d 75 6d 00 23 20 16|....supremum.#.|

0000c080 14 13 0c 06 00 00 10 0f 00 ba 00 00 00 2b 68 0b|.............+h.|

0000c090 00 00 00 00 00 06 53 80 00 00 00 32 01 10 61 62 62|.....S....2..abb|

0000c0a0 62 62 20 20 20 20 20 20 20 20 63 63 63 23 20 16|bb ccc#.|

0000c0b0 14 13 0c 06 00 00 18 0f 00 ea 00 00 00 2b 68 0c|.............+h.|

0000c0c0 00 00 00 00 06 53 80 00 00 00 32 01 1e 64 65 65|.....S....2..dee|

0000c0d0 65 65 20 20 20 20 20 20 20 20 66 66 66 21 9e 94|ee fff!..|

0000c0e0 14 13 0c 06 00 00 20 0f 00 74 00 00 00 2b 68 0d|........t...+h.|

0000c0f0 00 00 00 00 06 53 80 00 00 00 32 01 2c 64 00 00|.....S....2.,d..|

0000c100 00 00 00 00 00 00 00 00 66 66 66 00 00 00 00 00|........fff.....|
```

---

# 我们先来分析第一条记录：

---

```
23 20 16 14 13 0c 06/*变长字段长度列表*/
```

00 00 10 0f 00 ba/*Record Header，前面6个字节*/

00 00 00 2b 68 0b/*RowID*/

00 00 00 00 06 53/*TransactionID*/

80 00 00 00 32 01 10/*Roll Point*/

61/*第1列，'a'*/

62 62/*第2列，'bb'*/

62 62 20 20 20 20 20 20 20 20/*第3列，'bb'Char类型*/

63 63 63/*第4列，'ccc'*/

---

23 20 16 14 13 0c 06，看看06、0c、13、14、16、20、23，第一个字段长度为6，第二个字段为6，6+6=0x0C，第三个字段长度为7，6+6+7=0x13，第四个字段为1，6+6+7+1=0x14，第五个字段为2，6+6+7+1+2=0x16，第六个字段为10，6+6+7+1+2+10=0x20，第七个字段为3，6+6+7+1+2+10+3=0x23。

然后看看第二条记录，Record Header最后两个字节48，十进制22，32进制是0000000111，总共有7个字段，第三列字段长度为3，对应总长度33，第1列和第二列是隐藏的系统列。

我们把第四列的值清空，发现在Redundant行格式中，空值是不占用任何存储空间的，而最后一列的NULL值，依然占用存储空间。

---

21 9e 94 14 13 0c 06/*每个字段的长度*/

00 00 20 0f 00 74/*Record Header，前面6字节*/

00 00 00 2b 68 0d/*RowID*/

00 00 00 00 06 53/*TransactionID*/

80 00 00 00 32 01 10/*Roll Point*/

64/*第1列，'d'*/

00 00 00 00 00 00 00 00 00 00/*第3列，NULL*/

66 66 66/*第4列，'fff'*/

---

可以看到，Compact行格式下变长字段长度列表中存储的值依次是06 0c 13 14 94 9e 21，第4个变长字段（第5列NULL）对应的是94，而第6列CHAR列（非NULL）是9e，94+10=0x9e，第七列是21，就是14+3=0x21，这也说明了对于VARCHAR类型的NULL值，Redundant行格式是不占用存储空间的，而CHAR类型的NULL值是会占用的。

由于表mytest2使用的是Latin1字符集，所以表1中每个字符占用mytest2表中使用的是utf8字符集，CHAR类型在该字符集下占用10个字节，即10×3=30个字节。从Redundant行格式下对于CHAR类型会占用固定长度的存储空间这点，也可以看出上述的表述

# 4.3.3  VARCHAR类型

InnoDB存储引擎的表是索引组织表，这意味着数据是按主键组织存放的，而BLOB、LOB这样的大对象类型的数据不能放在同一个页上，因此实际存放这些BLOB类型数据的页需要单独申请，而在原有的页上保存VARCHAR类型数据的指针，这会使性能有所下降。

相比于VARCHAR，不同数据库对于DBA而言，MySQL数据库的VARCHAR类型可以存放比Oracle VARCHAR2最大的4000字节和SQL Server最大的8000字节还多的MySQL数据库VARCHAR类型最大长度为65535字节，但是实际能存放的字节数仍为65535字节，因为创建VARCHAR长度65535以下时，会报下面的错误：

---

```
mysql>CREATE TABLE test(

->a VARCHAR(65535)

->)CHARSET=latin1 ENGINE=InnoDB;

ERROR 1118(42000):Row size too large.The maximum row size for the used table type,not counting BLOBs,is 65535.You have to change
some columns to TEXT or BLOBs
```

---

从错误的消息中可以看出InnoDB存储引擎并不支持65535长度的VARCHAR类型，这是因为还有别的开销，通过测试发现能存放VARCHAR类型的最大长度为65532，因此创建下面的表就不会报错。

---

```
mysql>CREATE TABLE test(

->a VARCHAR(65532)

->)CHARSET=latin1 ENGINE=InnoDB;

Query OK,0 rows affected(0.15 sec)
```

---

需要注意上述的例子，在这个例子中的SQL_MODE限制这样的错误，如果更改这个参数，MySQL数据库会发出一个warning的提示。

---

```
mysql>CREATE TABLE test(

->a VARCHAR(65535)

->)CHARSET=latin1 ENGINE=InnoDB;
```

Query OK,0 rows affected,1 warning(0.14 sec)

mysql□SHOW WARNINGS\G;

***************************1.row***************************

Level:Note

Code:1246

Message:Converting column'a'from VARCHAR to TEXT

1 row in set(0.00 sec)

warning□□□□□□□□□□□□□□□MySQL□□□□□□□□VARCHAR□□□□□□□□
TEXT□□□□□test□□□□□□□□□□

mysql□SHOW CREATE TABLE test\G;

***************************1.row***************************

Table:test

Create Table:CREATE TABLE'test'(

'a'mediumtext

)ENGINE=InnoDB DEFAULT CHARSET=utf8

1 row in set(0.00 sec)

□□□□□□□□□□□□VARCHAR□□□□65 532□□□□□□□□□□□□latin1□□□□□□□GBK□
□UTF-8□□□□□□□□□□□□□□□□□

mysql□CREATE TABLE test(

-□a VARCHAR(65532)

-□)CHARSET=GBK ENGINE=InnoDB;

ERROR 1074(42000):Column length too big for column'a'(max=32767);use BLOB or TEXT instead

mysql□mysql□CREATE TABLE test(

-□a VARCHAR(65532)

-□)CHARSET=UTF8 ENGINE=InnoDB;

ERROR 1074(42000):Column length too big for column'a'(max=21845);use BLOB or TEXT instead

以此可以推断，VARCHAR最大为65532字节，因为还需要存放额外信息（max表示最大）。但实际上，规定行数据中任何VARCHAR（N）中的N，如果大于规定的最大字节数，则其限制和VARCHAR字段一样也为65535。注意，这里

还有一点，上述例子是在MySQL数据库下测得的65535长度。对于VARCHAR数据类型来说，其虽然可以存放最大长度为的数据，但是实际上却是

mysql□CREATE TABLE test2(

-□a VARCHAR(22000),

-□b VARCHAR(22000),

-□c VARCHAR(22000)

-□)CHARSET=latin1 ENGINE=InnoDB;

ERROR 1118(42000):Row size too large.The maximum row size for the used table type,not counting BLOBs,is 65535.You have to change some columns to TEXT or BLOBs

3个列相加已经超过了66000。但是InnoDB存储引擎同样不能创建字段长度总和大于65532的表。当然别忘了，InnoDB存储引擎的页为16KB，即16384字节，怎么能存放65532字节的数据呢？因此，这badInnoDB存储引擎同样不支持过长的字段，即B-tree node。所以这里还有一些的部分数据应放在称为Uncompress BLOB的页中来存放溢出的行数据。

mysql□CREATE TABLE t(

-□a VARCHAR(65532)

-□)ENGINE=InnoDB CHARSET=latin1;

Query OK,0 rows affected(0.15 sec)

mysql□INSERT INTO t

-□SELECT REPEAT('a',65532);

Query OK,1 row affected(0.08 sec)

Records:1 Duplicates:0 Warnings:0

上述例子创建了一个列，字段名为a，长度为65 532的VARCHAR表t，并插入了一个列a字段为65 532长度的数据。通过工具py_innodb_page_info，可以来分析表空间，可以看到如下图所示

```
[root@nineyou0-43 mytest]#py_innodb_page_info.py-v t.ibd

page offset 00000000,page type：File Space Header：

page offset 00000001,page type：Insert Buffer Bitmap：

page offset 00000002,page type：File Segment inode：

page offset 00000003,page type：B-tree Node：,page level：0000：

page offset 00000004,page type：Uncompressed BLOB Page：

page offset 00000005,page type：Uncompressed BLOB Page：

page offset 00000006,page type：Uncompressed BLOB Page：

page offset 00000007,page type：Uncompressed BLOB Page：

Total number of page:8:

Insert Buffer Bitmap:1

Uncompressed BLOB Page:4

File Space Header:1

B-tree Node:1

File Segment inode:1
```

从这个例子可以看到，由于我们保存的是一个列，一个B-tree Node，接着四个未压缩的二进制大对象页Uncompressed BLOB Page。在二进制大对象页中保存了之前插入的65532字节的数据。下面我们通过hexdump来进一步分析BLOB页中数据的存储方式。例如通过hexdump来查看偏移量为c000页的内容：

```
0000c000 67 ce fc 0b 00 00 00 03 ff ff ff ff ff ff ff ff|g..............|

0000c010 00 00 00 0a 6a d9 c0 89 45 bf 00 00 00 00 00 00|....j...E.......|

0000c020 00 00 00 00 00 c3 00 02 03 a7 80 03 00 00 00 00|...............|

0000c030 00 80 00 05 00 00 00 01 00 00 00 00 00 00 00 00|...............|

0000c040 00 00 00 00 00 00 00 00 01 a1 00 00 00 c3 00 00|...............|

0000c050 00 02 00 f2 00 00 00 c3 00 00 00 02 00 32 01 00|.............2..|

0000c060 02 00 1d 69 6e 66 69 6d 75 6d 00 02 00 0b 00 00|...infimum......|

0000c070 73 75 70 72 65 6d 75 6d 14 c3 00 00 00 10 ff f0|supremum........|

0000c080 00 00 00 b6 2b 00 00 00 00 51 4b 06 80 00 00 00|....+....QK.....|
```

```
0000c090 2d 01 10 61 61 61 61 61 61 61 61 61 61 61 61 61 61|-..aaaaaaaaaaaa|

0000c0a0 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61|aaaaaaaaaaaaaaaa|

0000c0b0 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61|aaaaaaaaaaaaaaaa|

0000c0c0 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61|aaaaaaaaaaaaaaaa|

0000c0d0 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61|aaaaaaaaaaaaaaaa|

0000c0e0 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61|aaaaaaaaaaaaaaaa|

0000c0f0 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61|aaaaaaaaaaaaaaaa|

0000c100 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61|aaaaaaaaaaaaaaaa|

0000c110 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61 61|aaaaaaaaaaaaaaaa|

......

0000c390 61 61 61 00 00 00 c3 00 00 00 04 00 00 00 26 00|aaa...........口.|

0000c3a0 00 00 00 00 00 00 fc fc 00 00 00 00 00 00 00 00 00|................|
```

从这可以看出0x0000c093到0x0000c392共768个字节为保存VARCHAR（65532）的前768个字节的prefix。由于前面保存的a字符的值都是相同的，因此只能通过最后的Uncompressed BLOB Page才能得到最终的结果，如图4-4所示。



图 4-4 行溢出数据的存储

如果存放VARCHAR数据类型的列值长度太长，使用了外部的BLOB页面来存放数据，而这些BLOB的页面类型是InnoDB，不属于任何索引，不管聚集还是辅助索引。因此B+Tree索引的高度都没有变化。需要注意的是，如果创建的是BLOB类型的列，那么在聚集或辅助索引中所保存的同样仅是指向真实数据的指针而已，InnoDB存储引擎并不关心BLOB的大小。从逻辑上来看，这些数据被视为"溢出页"，其实与外部存储的定义有异曲同工之妙。

---

```
mysql：CREATE TALBE t(

-：a VARCHAR(9000)

-：)ENGINE=InnoDB;

Query OK,0 rows affected(0.13 sec)

mysql：INSERT INTO t

-：SELECT REPEAT('a',9000);

Query OK,1 row affected(0.04 sec)

Records:1 Duplicates:0 Warnings:0

mysql：INSERT INTO t

-：SELECT REPEAT('a',9000);

Query OK,1 row affected(0.04 sec)

Records:1 Duplicates:0 Warnings:0
```

---

表t中插入了列a两个长度为9000的字符串，由于这时列的长度不能够完全放在一个页中（9000已经大于前面提到的阈值），我们再次使用工具py_innodb_page_info来观察，可以看到两个数据页和两个BLOB页。

---

```
[root@nineyou0-43 mytest]#py_innodb_page_info.py-v t.ibd

page offset 00000000,page type：File Space Header：

page offset 00000001,page type：Insert Buffer Bitmap：

page offset 00000002,page type：File Segment inode：

page offset 00000003,page type：B-tree Node：,page level：0000：

page offset 00000004,page type：Uncompressed BLOB Page：

page offset 00000005,page type：Uncompressed BLOB Page：

Total number of page:6:

Insert Buffer Bitmap:1
```

Uncompressed BLOB Page:2

File Space Header:1

B-tree Node:1

File Segment inode:1

---

## 接着 利用py_innodb_page_info工具来统计测试表的页类型分布，来分析验证行溢出数据的页类型

下面创建一个列为测试表，同时插入两条数据，列的类型为VARCHAR。在本例中可以看到行溢出数据是存放在BLOB页中的，创建表的列，同时插入的记录的列大小为8098，即创建了一个列类型为varchar，8098字节的行插入了2条记录。

---

```
mysql CREATE TABLE t(

- a varchar(8098)

)ENGINE=InnoDB;

Query OK,0 rows affected(0.12 sec)

mysql INSERT INTO t SELECT REPEAT('a',8098);

Query OK,1 row affected(0.04 sec)

Records:1 Duplicates:0 Warnings:0

mysql INSERT INTO t SELECT REPEAT('a',8098);

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0
```

---

## 通过py_innodb_page_info工具来分析t.ibd文件可以看到，数据存放在行溢出数据的页类型为BLOB页。对于熟悉Microsoft SQL Server数据库的DBA来说，可能对InnoDB存储引擎的VARCHAR类型有疑问，SQL Server中varchar（MAX）的存储

---

```
[root@nineyou0-43 mytest]#py_innodb_page_info.py-v t.ibd

page offset 00000000,page type File Space Header

page offset 00000001,page type Insert Buffer Bitmap
```

page offset 00000002,page type□File Segment inode□

page offset 00000003,page type□B-tree Node□,page level□0000□

page offset 00000000,page type□Freshly Allocated Page□

page offset 00000000,page type□Freshly Allocated Page□

Total number of page:6:

Freshly Allocated Page:2

Insert Buffer Bitmap:1

File Space Header:1

B-tree Node:1

File Segment inode:1

---

# □□□□□□□□□□TEXT□BLOB□□□□□□□□□□□□□□□□□□□□□Uncompressed BLOB Page□□□□□□□□□□□□□□□□□□□□□□□□□□BLOB□□□□□□□□□□□□VARCHAR□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□CREATE TABLE t(

-□a BLOB

-□)ENGINE=InnoDB;

Query OK,0 rows affected(0.12 sec)

mysql□INSERT INTO t SELECT REPEAT('a',8000);

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO t SELECT REPEAT('a',8000);

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO t SELECT REPEAT('a',8000);

Query OK,1 row affected(0.01 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO t SELECT REPEAT('a',8000);

Query OK,1 row affected(0.06 sec)

Records:1 Duplicates:0 Warnings:0

---

插入这条记录，因为BLOB没有超过行所能存放的最大长度8000，所以数据都存放在数据页中。通过工具py_innodb_page_info来查看表空间t.ibd，可以得到下面的结果，没有使用任何BLOB页：

```
[root@nineyou0-43 mytest]#py_innodb_page_info.py-v t.ibd

page offset 00000000,page type《File Space Header》

page offset 00000001,page type《Insert Buffer Bitmap》

page offset 00000002,page type《File Segment inode》

page offset 00000003,page type《B-tree Node》,page level《0001》

page offset 00000004,page type《B-tree Node》,page level《0000》

page offset 00000005,page type《B-tree Node》,page level《0000》

page offset 00000006,page type《B-tree Node》,page level《0000》

page offset 00000000,page type《Freshly Allocated Page》

Total number of page:8:

Freshly Allocated Page:1

Insert Buffer Bitmap:1

File Space Header:1

B-tree Node:4

File Segment inode:1
```

---

当然，如果可以预见到BLOB列的数据不会被经常读取到，那么将BLOB列的数据存放到外部的BLOB页中可能也是一种好的选择，这时可以强制BLOB列数据存放在溢出页中，保证数据页的768连续性。

## 4.3.4　Compressed和Dynamic行记录格式

InnoDB 1.0.x版本开始引入了新的文件格式（file format，可以理解为新的页格式），以前支持的Compact和Redundant格式称为Antelope文件格式，新的文件格式称为Barracuda文件格式。Barracuda文件格式下拥有两种新的行记录格式：Compressed和Dynamic。

新的两种记录格式对于存放在BLOB中的数据采用了完全的行溢出的方式，如图4-5所示，在数据页中只存放20个字节的指针，实际的数据都存放在Off Page中，而之前的Compact和Redundant两种格式会存放768个前缀字节。



图　4-5　Barracuda文件格式的行溢出

Compressed行记录格式的另一个功能就是，存储在其中的行数据会以zlib的算法进行压缩，因此对于BLOB、TEXT、VARCHAR这类大长度类型的数据能够进行非常有效的存储。

# 4.3.5　CHAR□□□□□□

□□□□VARCHAR□□□□□□□□□□□□□□CHAR□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□CHAR□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□latin1□□□□□MySQL 4.1□□□□□□CHR□N□□□□N□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□CHAR□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

```
mysql□CREATE TABLE j(

-□a CHAR(2)

-□)CHARSET=GBK ENGINE=InnoDB;

Query OK,0 rows affected(0.11 sec)

mysql□INSERT INTO j SELECT'ab';

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql□SET NAMES GBK;

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO j SELECT'□□';

Query OK,1 row affected(0.04 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO j SELECT'a';

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0
```

---

□□□□□□□□□□□□j□□□□□□□GBK□□□□□□□□□□□□□□□□□□□□□□'ab'□□'□□'□□□□□□□□□□□□□□□□□□□□□□□□

---

```
mysql□SELECT a,CHAR_LENGTH(a),LENGTH(a)

-□FROM j\G;
```

```
***************************1.row***************************

a:ab

CHAR_LENGTH(a):2

LENGTH(a):2

***************************2.row***************************

a:□□

CHAR_LENGTH(a):2

LENGTH(a):4

***************************3.row***************************

a:a

CHAR_LENGTH(a):1

LENGTH(a):1

3 rows in set(0.00 sec)
```

---

# 在这个例子中CHAR_LENGTH和CHAR字节计算的差异很明显，'ab'和'汉字'的字符数量都是2。但是字节数量'ab'就是2，而不是'汉字'就是4。下面再来看看HEX函数来计算字符的十六进制数值来加深理解。

---

```
mysql□SELECT a,HEX(a)

-□FROM j\G;

***************************1.row***************************

a:ab

HEX(a):6162

***************************2.row***************************

a:□□

HEX(a):CED2C3C7

***************************3.row***************************

a:a

HEX(a):61

3 rows in set(0.00 sec)
```

其它字符集下的结果。'ab'字段的内容就是0x6162，而汉字'中国'是0xCED2C3C7。这里顺便说明一下，CHAR字段类型存储时有一个小特性：如果不是采用UTF-8，CHAR（10）类型的字段并不总是占用10个字节。本来应该占用30个字节的，但这里只占用了两个字节，说明CHAR采用了变长。InnoDB在存储时做了优化，对于变长字符集，会将尾部的空格去掉。而CHAR字段类型具体的存储可以参考hexdump的结果，这里引用j.ibd文件：

0000c070 73 75 70 72 65 6d 75 6d 02 00 00 00 10 00 1c 00|supremum........|

0000c080 00 00 b6 2b 2b 00 00 00 51 52 da 80 00 00 00 2d|...++...QR.....-|

0000c090 01 10 61 62 04 00 00 00 18 ff d5 00 00 00 b6 2b|..ab...........+|

0000c0a0 2c 00 00 00 51 52 db 80 00 00 00 2d 01 10 ce d2|,...QR.....-....|

0000c0b0 c3 c7 00 00 00 00 00 00 00 00 00 00 00 00 00 00|................|

# 接着来分析这两条记录：

#第一条记录

02/*变长字段，2字节CHAR类型的长度*/

00/*NULL标志位*/

00 00 10 00 1c/*Recoder Header*/

00 00 00 b6 2b 2b/*RowID*/

00 00 00 51 52 da/*TransactionID*/

80 00 00 00 2d 01 10/*Roll Point*/

61 62/*列数'ab'*/

#第二条记录

04/*变长字段，4字节CHAR类型的长度*/

00/*NULL标志位*/

00 00 18 ff d5/*Recoder Header*/

00 00 00 b6 2b 2c/*RowID*/

00 00 00 51 52 db/*TransactionID*/

80 00 00 00 2d 01 10/*Roll Point*/

c3 d2 c3 c7/*汉字'测试'*/

#剩下的部分

02/*可变字段第2个（CHAR）的真实长度*/

00/*NULL标志位*/

00 00 20 ff b7/*Recoder Header*/

00 00 00 b6 2b 2d/*RowID*/

00 00 00 51 53 17/*TransactionID*/

80 00 00 00 2d 01 10/*Roll Point*/

61 20/*字段'a'*/

通过观察上面的记录，InnoDB其实并没有把CHAR类型的真实长度保存起来，因为CHAR类型是固定长度的，不够长度时在末尾填充空格，即0x20。InnoDB把填充空格后的字段值通过HEX字符集编码存储到聚集索引中。至此，我们已经说清楚了定长的CHAR和VARCHAR类型是如何存储在聚集索引上的。

# 4.4 InnoDB数据页结构

页是InnoDB存储引擎管理数据库的最小磁盘单位。页类型为InnoDB数据页的，其存放的即是表中行的实际数据了。本节将详细介绍InnoDB数据页的B-tree Node页的内部存储结构，希望通过本节的介绍能使读者更深入地理解InnoDB存储引擎页的内部结构。

注意　InnoDB数据页结构是本书新增的一个章节，当前在任何资料上都很难找到对于InnoDB数据页结构的详细解释。但是我觉得有必要了解InnoDB数据页的内部存储结构，因为Peter曾说过InnoDB存储引擎的Peter，他在《高性能MySQL》一书中，也只对于InnoDB存储引擎Compact行记录格式做了介绍，但对于页结构这部分，还没有任何的介绍。

InnoDB数据页由以下7个部分组成，如图4-6所示。

❑ File Header（文件头）

❑ Page Header（页头）

❑ Infimun和Supremum Records

❑ User Records（用户记录，即行记录）

❑ Free Space（空闲空间）

❑ Page Directory（页目录）

❑ File Trailer（文件结尾信息）

| File Header | 38字节 |
| Page Header | 56字节 |
| Infimun + Supremum Records | 行记录 |
| User Records | |
| Free Space | |
| Page Directory | |
| File Trailer | 8字节 |

Page Size

图 4-6　InnoDB数据页存储结构图

其中File Header、Page Header、File Trailer的大小是固定的，分别为38、56、8
字节，这些空间用来标记该页的一些信息，如Checksum、数据页所在B+树索引的层数等。User
Records、Free Space、Page Directory这些部分为实际的行记录存储空间，因此大小是动态
的。接下来部分将具体分析各个部分的存储内容。

## 4.4.1　File Header

File Header用来记录页的一些头信息，由表4-3的8个部分组成，共占用38字节。

表 4-3　File Header 组成部分

| 名　称 | 大小（字节） | 说　明 |
|---|---|---|
| FIL_PAGE_SPACE_OR_CHKSUM | 4 | 当 MySQL 为 MySQL4.0.14 之前的版本时，该值为 0。在之后的 MySQL 版本中，该值代表页的 checksum 值（一种新的 checksum 值） |
| FIL_PAGE_OFFSET | 4 | 表空间中页的偏移值。如某独立表空间 a.ibd 的大小为 1GB，如果页的大小为 16KB，那么总共有 65 536 个页。FIL_PAGE_OFFSET 表示该页在所有页中的位置。若此表空间的 ID 为 10，那么搜索页（10，1）就表示查找表 a 中的第二个页 |
| FIL_PAGE_PREV | 4 | 当前页的上一个页，B+ Tree 特性决定了叶子节点必须是双向列表 |
| FIL_PAGE_NEXT | 4 | 当前页的下一个页，B+ Tree 特性决定了叶子节点必须是双向列表 |
| FIL_PAGE_LSN | 8 | 该值代表该页最后被修改的日志序列位置 LSN（Log Sequence Number） |
| FIL_PAGE_TYPE | 2 | InnoDB 存储引擎页的类型。常见的类型见表 4-4。记住 0x45BF，该值代表了存放的是数据页，即实际行记录的存储空间 |
| FIL_PAGE_FILE_FLUSH_LSN | 8 | 该值仅在系统表空间的一个页中定义，代表文件至少被更新到了该 LSN 值。对于独立表空间，该值都为 0 |
| FIL_PAGE_ARCH_LOG_NO_OR_SPACE_ID | 4 | 从 MySQL 4.1 开始，该值代表页属于哪个表空间 |

表 4-4 InnoDB 存储引擎中页的类型

| 名 称 | 十六进制 | 解 释 |
|---|---|---|
| FIL_PAGE_INDEX | 0x45BF | B+ 树叶节点 |
| FIL_PAGE_UNDO_LOG | 0x0002 | Undo Log 页 |
| FIL_PAGE_INODE | 0x0003 | 索引节点 |
| FIL_PAGE_IBUF_FREE_LIST | 0x0004 | Insert Buffer 空闲列表 |
| FIL_PAGE_TYPE_ALLOCATED | 0x0000 | 该页为最新分配 |
| FIL_PAGE_IBUF_BITMAP | 0x0005 | Insert Buffer 位图 |
| FIL_PAGE_TYPE_SYS | 0x0006 | 系统页 |
| FIL_PAGE_TYPE_TRX_SYS | 0x0007 | 事务系统数据 |
| FIL_PAGE_TYPE_FSP_HDR | 0x0008 | File Space Header |
| FIL_PAGE_TYPE_XDES | 0x0009 | 扩展描述页 |
| FIL_PAGE_TYPE_BLOB | 0x000A | BLOB 页 |

## 4.4.2　Page Header

紧跟File Header之后的是Page Header，它占用固定的字节数，其长度为14，其所在的页偏移量为56。如表格4-5所示。

表 4-5  Page Header 组成部分

| 名　称 | 大小（字节） | 说　明 |
|---|---|---|
| PAGE_N_DIR_SLOTS | 2 | 在 Page Directory（页目录）中的 Slot（槽）数，"4.4.5 Page Directory"小节中会介绍 |
| PAGE_HEAP_TOP | 2 | 堆中第一个记录的指针，记录在页中是根据堆的形式存放的 |
| PAGE_N_HEAP | 2 | 堆中的记录数。一共占用 2 字节，但是第 15 位表示行记录格式 |
| PAGE_FREE | 2 | 指向可重用空间的首指针 |
| PAGE_GARBAGE | 2 | 已删除记录的字节数，即行记录结构中 delete flag 为 1 的记录大小的总数 |
| PAGE_LAST_INSERT | 2 | 最后插入记录的位置 |
| PAGE_DIRECTION | 2 | 最后插入的方向。可能的取值为：<br>❏ PAGE_LEFT（0x01）<br>❏ PAGE_RIGHT（0x02）<br>❏ PAGE_SAME_REC（0x03）<br>❏ PAGE_SAME_PAGE（0x04）<br>❏ PAGE_NO_DIRECTION（0x05） |
| PAGE_N_DIRECTION | 2 | 一个方向连续插入记录的数量 |
| PAGE_N_RECS | 2 | 该页中记录的数量 |
| PAGE_MAX_TRX_ID | 8 | 修改当前页的最大事务 ID，注意该值仅在 Secondary Index 中定义 |
| PAGE_LEVEL | 2 | 当前页在索引树中的位置，0x00 代表叶节点，即叶节点总是在第 0 层 |
| PAGE_INDEX_ID | 8 | 索引 ID，表示当前页属于哪个索引 |

| 名　称 | 大小（字节） | 说　明 |
|---|---|---|
| PAGE_BTR_SEG_LEAF | 10 | B+ 树数据页非叶节点所在段的 segment header。注意该值仅在 B+ 树的 Root 页中定义 |
| PAGE_BTR_SEG_TOP | 10 | B+ 树数据页所在段的 segment header。注意该值仅在 B+ 树的 Root 页中定义 |

## 4.4.3　Infimum、Supremum Record

在InnoDB数据页中，对于用户的数据记录，总是存在着两条伪记录，分别是Infimum（下确界）和伪记录是最小值。另外一条是Supremum（上确界）伪记录是最大值。这两条伪记录是在数据页创建的时候被创建的，并且在任何情况下都不会被删除。在Compact格式和Redundant格式中，两条伪记录的存储方式也是不同的。图4-7显示了Infimum、Supremum的结构。

B+ Tree

页

| Infinum Record | 行记录数据 | 行记录数据 | 行记录数据 | ... | Supremum Record |

图 4-7 Infinum、Supremum Record

## 4.4.4　User Record、Free Space

User Record中存储的是实际的用户记录数据，这些记录通过InnoDB提供的索引组织成B+树结构进行管理。

Free Space就是空闲空间，新插入的记录如果没有可以重用的已删除记录空间，就会从这里进行分配。

## 4.4.5　Page Directory

Page Directory□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Slots□□□□□□□□□□Directory Slots□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□sparse directory□□□□□□□□□□□□□□□□□□□□□□□□Infimum□n_owned□□□□□1□□□□Supremum□n_owned□□□□□□□[1□8]□□□□□□□□□n_owned□□□□□□□[4□8]□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□Slots□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□'i'□'d'□'c'□'b'□'e'□'g'□'l'□'h'□'f'□'j'□'k'□'a'□□□□□□□□□□□□□4□□□□□Slots□□□□□□□□□'a'□'e'□'i'□□

□□□□InnoDB□□□□□□□Page Direcotry□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□recorder header□□next_record□□□□□□□□□□□□□□Page Directory□□□□□□recorder header□□n_owned□□□□□□□□□□□□□□□□□□□□Page Directory□□□

□□□□□□□□□B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Page Directory□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 4.4.6　File Trailer

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□File Trailer□□□

File Trailer□□□□FIL_PAGE_END_LSN□□□□□8□□□□□4□□□□□□□□checksum□□□□□4□□□File Header□□FIL_PAGE_LSN□□□□□□□□□□□File Header□□FIL_PAGE_SPACE_OR_CHKSUM□FIL_PAGE_LSN□□□□□□□□□□□□□checksum□□□□□□□□InnoDB□checksum□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□not corrupted□□

□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Corrupt□□□□□□□File Trailer□□□□□□□□□□□□□□□□□□□□□□□□□□□□innodb_checksums□□□□□□□□□□□□□□□□□□

MySQL 5.6.6□□□□□□□□□□innodb_checksum_algorithm□□□□□□□□□□checksum□□□□□□□□□□crc32□□□□□□□□innodb□crc32□none□strict_innodb□strict_crc32□strict_none□

innodb□□□□□□□□InnoDB□□checksum□□□□□□crc32□MySQL 5.6.6□□□□□□checksum□□□□□□□□innodb□□□□□□□□□□□□□□□□checksum□□□strict□□□□□□□□□□□□MySQL□□□□□□□□□□□□none□□□□□□□□checksum□□□□

strict_*□□□□□□□□□□□□□□□checksum□□□□□□□□□□□□□□□MySQL□□□□□□□MySQL 5.6.6□□□□□□□□□□strict_crc32□□□□□□□□□□□□□□strict_crc32□□□□□□□□□□□□□□□□innodb□crc32□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□mysql_upgrade□□□□

# 4.4.7 InnoDB自增长与字段自增长

到这里需要说明的是，对于序列，InnoDB有着和其他数据库完全不同的策略来实现自增长的功能。为了具体说明这个问题，先来看下面的一个例子，在例子中我们创建一个表t，其中列是自增长的。

---

mysql＞DROP TABLE IF EXISTS t;

Query OK,0 rows affected(0.04 sec)

mysql＞CREATE TABLE t(

-＞a INT UNSIGNED NOT NULL AUTO_INCREMENT,

-＞b CHAR(10),

-＞PRIMARY KEY(a),

-＞)ENGINE=InnoDB CHARSET=UTF8;

Query OK,0 rows affected(0.00 sec)

mysql＞DELIMITER$$

mysql＞CREATE PROCEDURE load_t(count INT UNSIGNED)

-＞BEGIN

-＞SET@c=0;

-＞WHILE@c＜count DO

-＞INSERT INTO t

-＞SELECT NULL,REPEAT(CHAR(97+RAND()*26),10);

-＞SET@c=@c+1;

-＞END WHILE;

-＞END;

-＞$$

Query OK,0 rows affected(0.00 sec)

mysql＞DELIMITER;

mysql＞CALL load_t(100);

Query OK,0 rows affected(0.60 sec)

```
mysql SELECTa,bFROM t LIMIT 10;

+----+--------------+

|a|b|

+----+--------------+

|1|dddddddddd|

|2|hhhhhhhhhh|

|3|bbbbbbbbbb|

|4|iiiiiiiiii|

|5|nnnnnnnnnn|

|6|qqqqqqqqqq|

|7|oooooooooo|

|8|yyyyyyyyyy|

|9|yyyyyyyyyy|

|10|vvvvvvvvvv|

+----+--------------+

10 rows in set(0.00 sec)
```

## 这时再用py_innodb_page_info工具查看t.ibd，可以看到如下内容:

```
[root@nineyou0-43 mytest]#py_innodb_page_info.py-v t.ibd

page offset 00000000,page type File Space Header

page offset 00000001,page type Insert Buffer Bitmap

page offset 00000002,page type File Segment inode

page offset 00000003,page type B-tree Node ,page level 0000

page offset 00000000,page type Freshly Allocated Page

page offset 00000000,page type Freshly Allocated Page

Total number of page:6:

Freshly Allocated Page:2

Insert Buffer Bitmap:1

File Space Header:1
```

B-tree Node:1

File Segment inode:1

---

# 为了进一步证明page offset 3就是我们需要的，hexdump整个t.ibd，搜索其中的内容，发现偏移量是0x0000c000（16K*3=0xc000）的数据页就是我们要找的数据页

---

```
0000c000 52 1b 24 00 00 00 00 00 03 ff ff ff ff ff ff ff|R.$.............|
0000c010 00 00 00 0a 6a e0 ac 93 45 bf 00 00 00 00 00 00|....j...E.......|
0000c020 00 00 00 00 00 dc 00 1a 0d c0 80 66 00 00 00 00|...........f....|
0000c030 0d a5 00 02 00 63 00 64 00 00 00 00 00 00 00 00|.....c.d........|
0000c040 00 00 00 00 00 00 00 00 01 ba 00 00 00 dc 00 00|................|
0000c050 00 02 00 f2 00 00 00 dc 00 00 00 02 00 32 01 00|.............2..|
0000c060 02 00 1c 69 6e 66 69 6d 75 6d 00 05 00 0b 00 00|...infimum......|
0000c070 73 75 70 72 65 6d 75 6d 0a 00 00 00 10 00 22 00|supremum......".|
0000c080 00 00 01 00 00 00 51 6d eb 80 00 00 00 2d 01 10|......Qm.....-..|
0000c090 64 64 64 64 64 64 64 64 64 64 0a 00 00 00 18 00|dddddddddd......|
0000c0a0 22 00 00 00 02 00 00 00 51 6d ec 80 00 00 00 2d|".......Qm.....-|
0000c0b0 01 10 68 68 68 68 68 68 68 68 68 68 0a 00 00 00|..hhhhhhhhhh....|
0000c0c0 20 00 22 00 00 00 03 00 00 00 51 6d ed 80 00 00| .".......Qm....|
0000c0d0 00 2d 01 10 62 62 62 62 62 62 62 62 62 62 0a 00|.-..bbbbbbbbbb..|
0000c0e0 04 00 28 00 22 00 00 00 04 00 00 00 51 6d ee 80|..(.".......Qm..|
0000c0f0 00 00 00 2d 01 10 69 69 69 69 69 69 69 69 69 69|...-..iiiiiiiiii|
0000c100 0a 00 00 00 30 00 22 00 00 00 05 00 00 00 51 6d|....0.".......Qm|
0000c110 ef 80 00 00 00 2d 01 10 6e 6e 6e 6e 6e 6e 6e 6e|.....-..nnnnnnnn|
0000c120 6e 6e 0a 00 00 00 38 00 22 00 00 00 06 00 00 00|nn....8.".......|
0000c130 51 6d f0 80 00 00 00 2d 01 10 71 71 71 71 71 71|Qm.....-..qqqqqq|
0000c140 71 71 71 71 0a 00 00 00 40 00 22 00 00 00 07 00|qqqq....@.".....|
0000c150 00 00 51 6d f1 80 00 00 00 2d 01 10 6f 6f 6f 6f|..Qm.....-..oooo|
0000c160 6f 6f 6f 6f 6f 6f 0a 00 04 00 48 00 22 00 00 00|oooooo....H."...|
```

```
0000c170 08 00 00 00 51 6d f2 80 00 00 00 2d 01 10 79 79|....Qm.....-..yy|

0000c180 79 79 79 79 79 79 79 79 0a 00 00 00 50 00 22 00|yyyyyyyy....P.".|

0000c190 00 00 09 00 00 00 51 6d f3 80 00 00 00 2d 01 10|......Qm.....-..|

0000c1a0 79 79 79 79 79 79 79 79 79 79 0a 00 00 00 58 00|yyyyyyyyyy....X.|

0000c1b0 22 00 00 00 0a 00 00 00 51 6d f4 80 00 00 00 2d|".......Qm.....-|

0000c1c0 01 10 76 76 76 76 76 76 76 76 76 76 0a 00 00 00|..vvvvvvvvvv....|

0000c1d0 60 00 22 00 00 00 0b 00 00 00 51 6d f5 80 00 00|'.".......Qm....|

0000c1e0 00 2d 01 10 6b 6b 6b 6b 6b 6b 6b 6b 6b 6b 0a 00|.-..kkkkkkkkkk..|

0000c1f0 04 00 68 00 22 00 00 00 0c 00 00 00 51 6d f6 80|..h.".......Qm..|

……

0000ffc0 00 00 00 00 00 00 70 0d 1d 0c 95 0c 0d 0b 85 0a fd|.....p..........|

0000ffd0 0a 75 09 ed 09 65 08 dd 08 55 07 cd 07 45 06 bd|.u...e...U...E..|

0000ffe0 06 35 05 ad 05 25 04 9d 04 15 03 8d 03 05 02 7d|.5...%.........}|

0000fff0 01 f5 01 6d 00 e5 00 63 95 ae 5d 39 6a e0 ac 93|...m...c..]9j...|
```

# 最后来看看File Header（38字节）

· 52 1b 24 00（页面的的Checksum）；

· 00 00 00 03（页号，这里是0号页）；

· ff ff ff ff（同一个区的上一个页的页号，这里没有上一个页，0xffffffff）；

· ff ff ff ff（同一个区的下一个页的页号，这里没有下一个页，0xffffffff）；

· 00 00 00 0a 6a e0 ac 93（页的LSN）；

· 45 bf（页的类型，0x45bf，数据页）；

· 00 00 00 00 00 00 00（文件至少被刷新到的）；

· 00 00 00 dc（页属于的SPACE ID）；

最后，我们看看Page Header，由于它紧跟着File Trailer之后，而File Trailer之前是File Header，所以我们先看看紧跟在File Trailer后8个字节：

---

95 ae 5d 39 6a e0 ac 93

---

□95 ae 5d 39（Checksum）：这也是个checksum，它与File Header开头的checksum算法是一样的。

□6a e0 ac 93：它是页面File Header中存储的LSN的4个低位字节

接下来的56字节是Page Header，我们先看看页面的Page Header里存储的信息，然后再结合数据进行分析：

---

Page Header(56 bytes):

PAGE_N_DIR_SLOTS=0x001a

PAGE_HEAP_TOP=0x0dc0

PAGE_N_HEAP=0x8066

PAGE_FREE=0x0000

PAGE_GARBAGE=0x0000

PAGE_LAST_INSERT=0x0da5

PAGE_DIRECTION=0x0002

PAGE_N_DIRECTION=0x0063

PAGE_N_RECS=0x0064

PAGE_MAX_TRX_ID=0x0000000000000000

PAGE_LEVEL=00 00

PAGE_INDEX_ID=0x00000000000001ba

PAGE_BTR_SEG_LEAF=0x000000dc0000000200f2

PAGE_BTR_SEG_TOP=0x000000dc000000020032

---

PAGE_N_DIR_SLOTS=0x001a：表示Page Directory有26个槽，每个槽占用2个字节，也就是从0x0000ffc4到0x0000fff7这段空间代表的

```
0000ffc0 00 00 00 00 00 70 0d 1d 0c 95 0c 0d 0b 85 0a fd|.....p..........|

0000ffd0 0a 75 09 ed 09 65 08 dd 08 55 07 cd 07 45 06 bd|.u...e...U...E..|

0000ffe0 06 35 05 ad 05 25 04 9d 04 15 03 8d 03 05 02 7d|.5...%.........}|

0000fff0 01 f5 01 6d 00 e5 00 63 95 ae 5d 39 6a e0 ac 93|...m...c..]9j...|
```

# PAGE_HEAP_TOP=0x0dc0，根据前文可知，该数据页地址为 0xc000+0x0dc0=0xcdc0，这个偏移的空间已经被使用，这个位置之后的就是可利用的剩余空间了。

```
0000cdb0 00 00 00 2d 01 10 70 70 70 70 70 70 70 70 70 70|...-..pppppppppp|

0000cdc0 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00|................|

0000cdd0 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00|................|

0000cde0 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00|................|
```

# PAGE_N_HEAP=0x8066，表示该页格式为 Compact，并且总共有 0x0802，即两条（Redundant 格式固定为 2）虚拟记录，也就是我们之前说的 Infinimun、Supremum。该页中共有 0x8066-0x8002=0x64，也就是十进制的 100 条记录。

# PAGE_FREE=0x0000，表示被删除的记录链表的头指针，因为目前没有被删除的记录，所以值为 0。

# PAGE_GARBAGE=0x0000，表示被删除的字节为 0。因为目前没有被删除的记录，所以值为 0。

# PAGE_LAST_INSERT=0x0da5，表示最后插入记录的位置，也就是这个数据页地址 0xc0000+0x0da5=0xcda5，就是下面的：

```
0000cda0 00 03 28 f2 cb 00 00 00 64 00 00 00 51 6e 4e 80|..(.....d...QnN.|

0000cdb0 00 00 00 2d 01 10 70 70 70 70 70 70 70 70 70 70|...-..pppppppppp|

0000cdc0 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00|................|
```

上，所有记录都是按照a列从100开始递增插入的，所以接下来的插入的方向都是向右，所以其值为向右。

PAGE_DIRECTION=0x0002，因为我们是按照主键递增的顺序插入的，所以PAGE_DIRECTION的方向向右，为0x00002。

PAGE_N_DIRECTION=0x0063，表示一个方向连续插入记录的数量，因为我们已经插入了100条，去掉第一条，为99。

PAGE_N_RECS=0x0064，本页面中的记录数量为100。紧接着的是PAGE_N_HEAP，我们的PAGE_N_HEAP这两个字节这个记录，与数据的记录数量不一致，为0x8066。

PAGE_LEVEL=0x00，当前页面在索引中的层级，叶子节点的层级为，从下往上的B+数的层级，因为B+树很矮，所以为0x00。

PAGE_INDEX_ID=0x00000000000001ba，索引ID。

其他的属性我们对Page Header的各个属性都已经了解的差不多了，接下来我们就看InnoDB存储引擎两个非常特殊的记录，最小记录和最大记录。

---

0000c050 00 02 00 f2 00 00 00 dc 00 00 00 02 00 32 01 00|.............2..|

0000c060 02 00 1c 69 6e 66 69 6d 75 6d 00 05 00 0b 00 00|...infimum......|

0000c070 73 75 70 72 65 6d 75 6d 0a 00 00 00 10 00 22 00|supremum......".|

---

我们0xc05E到0xc077这一块区域是两个特殊的记录，因为InnoDB存储引擎的行记录格式为的格式，Char（8个字节），我们接下来就按照行记录的格式对两个特殊的记录进行分析。

---

#Infimum记录信息

01 00 02 00 1c/*recorder header*/

69 6e 66 69 6d 75 6d 00/*这几个字节表示存储的内容，为Infimum这个单词。

/*0x00结束符*/

#Supremum记录信息

05 00 0b 00 00/*recorder header*/

73 75 70 72 65 6d 75 6d/*这个不管它是什么，因为它是Supremum*/

---

## 接下来就是infimum记录的recorder header信息，最后两个字节是00 1c，代表的是第 一条用户记录的相对位置，也就是0xc063+0x001c就是0xc07f，0xc07f就是 下面这条用户数据信息，看看是不是我们想要的那条数据。

---

0000c070 73 75 70 72 65 6d 75 6d 0a 00 00 00 10 00 22 00|supremum......".|

0000c080 00 00 01 00 00 00 51 6d eb 80 00 00 00 2d 01 10|......Qm.....-..|

0000c090 64 64 64 64 64 64 64 64 64 64 0a 00 00 00 18 00|dddddddddd......|

0000c0a0 22 00 00 00 02 00 00 00 51 6d ec 80 00 00 00 2d|".......Qm.....-|

---

## 我们看看下面这条数据是不是我们要查询的，拆分一下这条用户数据记录。

---

/*用户数据记录*/

00 00 00 01/*因为我们没有设置主键，所以有ROWID，也就a的值1*/

00 00 00 51 6d eb/*Transaction ID*/

80 00 00 00 2d 01 10/*Roll Pointer*/

64 64 64 64 64 64 64 64 64 64/*b的内容'aaaaaaaaaa'*/

---

## 我们用命令查询一下，看看结果。

---

mysql》SELECT a,b,hex(b)FROM t ORDER BY a LIMIT 1;

+---+------------+-------------------------------+

|a|b|hex(b)|

+---+------------+-------------------------------+

|1|dddddddddd|64646464646464646464|

+---+------------+-------------------------------+

1 row in set(0.00 sec)

---

最后Recorder Header中保存着该记录与下一记录的相对位置偏移量，我们可以利用这个相对位置偏移量找到下一个记录，再结合Page Header的PAGE_PREV、PAGE_NEXT等字段，我们就可以完整复现InnoDB中整个逻辑链表的串联与访问逻辑。

接下来是Page Directory，其起始偏移量为0x0000ffc4，0x0000fff7，我们把Page Directory打印出来

---

```
0000ffc0 00 00 00 00 00 70 0d 1d 0c 95 0c 0d 0b 85 0a fd|.....p..........|

0000ffd0 0a 75 09 ed 09 65 08 dd 08 55 07 cd 07 45 06 bd|.u...e...U...E..|

0000ffe0 06 35 05 ad 05 25 04 9d 04 15 03 8d 03 05 02 7d|.5...%.........}|

0000fff0 01 f5 01 6d 00 e5 00 63 95 ae 5d 39 6a e0 ac 93|...m...c..]9j...|
```

---

我们以表格开头Page Directory中的第一个槽位为例，其前2个字节的数据为00 63，即相对偏移量为0xc063，00 70表示第二个槽位的相对偏移量为0xc070，这两个槽位分别对应Infimum、Supremum，这表明Page Directory中的槽位是从记录的逻辑顺序的尾部向头部逆序排列的，这种设计方便InnoDB进行二分查找。我们结合Recorder Header的n_owned字段进行分析，InnoDB规定每个槽位最多管理a个5条记录，我们继续看Page Directory中的第三个槽位，其数据为00 e5，即相对偏移量为0xc0e5。

---

```
0000c0e0 04 00 28 00 22 00 00 00 04 00 00 00 51 6d ee 80|..(."......Qm..|

0000c0f0 00 00 00 2d 01 10 69 69 69 69 69 69 69 69 69 69|...-..iiiiiiiiii|

0000c100 0a 00 00 00 30 00 22 00 00 00 05 00 00 00 51 6d|....0."......Qm|

0000c110 ef 80 00 00 00 2d 01 10 6e 6e 6e 6e 6e 6e 6e 6e|.....-..nnnnnnnn|

0000c120 6e 6e 0a 00 00 00 38 00 22 00 00 00 06 00 00 00|nn....8.".......|

0000c130 51 6d f0 80 00 00 00 2d 01 10 71 71 71 71 71 71|Qm.....-..qqqqqq|

0000c140 71 71 71 71 0a 00 00 00 40 00 22 00 00 00 07 00|qqqq....@."....|
```

---

该槽位指向记录的第4条记录，我们找到第6条记录，其前面的5字节是Record Header，04 00 28 00 22，其中4、8字节为n_owned，这表明该槽位管理4条记录，从第4条记录开始，到该记录截止。我们看recorder header的后两个字节，其值为0x0022，即下一个记录的偏移量为0xc107，即在逻辑链表中第5条记录。

把线下整理的树传上去。从这个过程可以看出，还原逻辑备份对线上服务的影响非常大。最后，我们再看看InnoDB是怎样简单实现增删改查的。

# 4.5 Named File Formats□□

随着InnoDB存储引擎的发展，新的文件格式被用来支持新的功能。例如之前我们介绍过InnoDB 1.0.x版本提供了新的页类型（新的文件格式）来支持压缩功能、Off page功能（新的行溢出存储方式）。如果要正确识别新老文件格式的不同，就需要InnoDB 1.0.x版本引入InnoDB存储引擎的Named File Formats机制来解决不同版本下页结构兼容性的问题。

InnoDB存储引擎在1.0.x版本之前的文件格式（file format）都定义为Antelope，这个版本支持之前版本的页结构。新的文件格式命名为Barracuda，其支持之前版本的页结构。新的文件格式命名如图4-8所示。Barracuda文件格式和Antelope文件格式对比，Antelope文件格式有Compact和Redudant两种，而Barracuda文件格式又包括Antelope所有的文件格式，此外新增了两种页结构，即Compressed和Dynamic两种格式。



图 4-8　文件格式

InnoDB Pluginでは、ファイルフォーマットはInnoDBの重要機能を利用する際の予約語のよう
な役割となっています。また、Appleが動物の名前を利用しているのと同様に、ファイルフォー
マットには動物の名前が利用されます。

---

/**List of animal names representing file format.*/

```c
static const char*file_format_name_map[]={

"Antelope",

"Barracuda",

"Cheetah",

"Dragon",

"Elk",

"Fox",

"Gazelle",

"Hornet",

"Impala",

"Jaguar",

"Kangaroo",

"Leopard",

"Moose",

"Nautilus",

"Ocelot",

"Porpoise",

"Quail",

"Rabbit",

"Shark",

"Tiger",

"Urchin",

"Viper",

"Whale",

"Xenops",
```

"Yak",

"Zebra"

};

---

## 其中innodb_file_format这个参数以前是没有的，这个参数只有在新版本的InnoDB存储引擎中才有。

---

mysql（SELECT@@version\G;

***************************1.row***************************

@@version:5.1.37

1 row in set(0.00 sec)

mysql（SHOW VARIABLES LIKE'innodb_version'\G;

***************************1.row***************************

Variable_name:innodb_version

Value:1.0.4

1 row in set(0.00 sec)

mysql（SHOW VARIABLES LIKE'innodb_file_format'\G;

***************************1.row***************************

Variable_name:innodb_file_format

Value:Barracuda

1 row in set(0.00 sec)

---

## 参数innodb_file_format_check用来检测当前InnoDB存储引擎文件格式的支持度，该值默认为ON，如果出现不支持的情况，可能会出现如下的错误提示：

---

InnoDB:Warning:the system tablespace is in a

file format that this version doesn't support

---

# 4.6 约束

## 4.6.1 约束的类型

数据完整性是指数据的可靠性和准确性，数据完整性可以通过完整性约束来保证，而完整性约束是保证数据库中数据完整性的制度。换句话说，完整性约束（constraint）是关系数据库的一个重要概念，它用来保证数据库中数据的完整性和正确性。

关系型数据库系统和文件系统的一个不同点是，InnoDB存储引擎本身提供了大量数据完整性检验。例如Primary Key、Unique Key等的约束，这些约束是在表创建时定义的，用来保证数据的完整性。

一般来说，数据完整性有以下三种形式，InnoDB存储引擎提供以下几种约束来保证实现这三种完整性。

□实体完整性，保证表中有一个主键，不可重复。

□通过□Foreign Key，实现。

□域完整性。

□通过以下□DEFAULT，约束实现，实现数据的有效性。

□参照完整性，保证两张表之间的InnoDB存储引擎提供了外键约束来实现，保证数据的有效性，关系的。

对于InnoDB存储引擎本身而言，提供了以下约束：

□Primary Key

□Unique Key

□Foreign Key

□Default

□NOT NULL

# 4.6.2 □□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□

□□□ALTER TABLE□□□□□□□□□□

□Unique Key□□□□□□□□□□□□□□□□□□□□□□CREATE UNIQUE INDEX□□□□□□□□□□□□□□□□□□□□□□□□□□PRIMARY□□□□□Unique Key□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Unique Key□□□□□□□□Foreign Key□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

```
mysql□CREATE TABLE u(

-□id INT,

-□name VARCHAR(20),

-□id_card CHAR(18),

-□PRIMARY KEY(id),

-□UNIQUE KEY(name));

Query OK,0 rows affected(0.16 sec)

mysql□SELECT constraint_name,constraint_type

-□FROM

-□information_schema.TABLE_CONSTRAINTS

-□WHERE table_schema='mytest'AND table_name='u';\G;

***************************1.row***************************

constraint_name:PRIMARY

constraint_type:PRIMARY KEY

***************************2.row***************************

constraint_name:name

constraint_type:UNIQUE

2 rows in set(0.00 sec)
```

有一个唯一性约束，它的名字就是PRIMARY。下面向其中增加一个唯一性约束，然后再来查看。使用的是ALTER TABLE语句，这个语句将在后面的章节中讲到。先来看这个例子:

mysql〉ALTER TABLE u

-〉ADD UNIQUE KEY uk_id_card(id_card);

Query OK,0 rows affected(0.19 sec)

Records:0 Duplicates:0 Warnings:0

mysql〉SELECT constraint_name,constraint_type

-〉FROM

-〉information_schema.TABLE_CONSTRAINTS

-〉WHERE table_schema='mytest'AND table_name='u';\G;

***************************1.row***************************

constraint_name:PRIMARY

constraint_type:PRIMARY KEY

***************************2.row***************************

constraint_name:name

constraint_type:UNIQUE

***************************3.row***************************

constraint_name:uk_id_card

constraint_type:UNIQUE

3 rows in set(0.00 sec)

## 外键（Foreign Key）约束。下面创建一个Foreign Key约束，首先要创建一个被引用的表，这里引用表p。

mysql〉CREATE TABLE p(

-〉id INT,

-〉u_id INT,

-〉PRIMARY KEY(id),

-　FOREIGN KEY(u_id)REFERENCES p(id));

Query OK,0 rows affected(0.13 sec)

mysql　SELECT constraint_name,constraint_type

-　FROM

-　information_schema.TABLE_CONSTRAINTS

-　WHERE table_schema='mytest'and table_name='p'\G;

***************************1.row***************************

constraint_name:PRIMARY

constraint_type:PRIMARY KEY

***************************2.row***************************

constraint_name:p_ibfk_1

constraint_type:FOREIGN KEY

2 rows in set(0.00 sec)

---

# 　　　　　　　　information_schema　　　　TABLE_CONSTRAINTS　　　　MySQL　　　　　　　　　　Foreign Key　　　　　　　　　　　　　　REFERENTIAL_CONSTRAINTS　　　　　　　　　　　　　　　

---

mysql　SELECT*FROM

-　information_schema.REFERENTIAL_CONSTRAINTS

-　WHERE constraint_schema='mytest'\G;

***************************1.row***************************

CONSTRAINT_CATALOG:NULL

CONSTRAINT_SCHEMA:test2

CONSTRAINT_NAME:p_ibfk_1

UNIQUE_CONSTRAINT_CATALOG:NULL

UNIQUE_CONSTRAINT_SCHEMA:test2

UNIQUE_CONSTRAINT_NAME:PRIMARY

MATCH_OPTION:NONE

UPDATE_RULE:RESTRICT

DELETE_RULE:RESTRICT

TABLE_NAME:p

REFERENCED_TABLE_NAME:p

1 row in set(0.00 sec)

# 4.6.3 　索引与约束的关系

索引的作用是加速查询，Primary Key、Unique Key这种约束既是约束，又是索引，它们既要求某列的数据满足一定的约束条件，又能加速查询。

与这些约束相对应的是普通索引，普通索引的作用仅仅是加速查询，它本身对插入表中的数据没有任何约束，只是在表中增加普通索引以后，数据库系统会为普通索引建立相应的数据结构以加速查询的效率。

# 4.6.4 □□□□□□□□□

□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□NOT NULL□□□□□□□□□NULL□□MySQL□□□□□□□□□□0□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□CREATE TABLE a(

-□id INT NOT NULL,

-□date DATE NOT NULL);

Query OK,0 rows affected(0.13 sec)

mysql□INSERT INTO a

-□SELECT NULL,'2009-02-30';

Query OK,1 row affected,2 warnings(0.04 sec)

Records:1 Duplicates:0 Warnings:2

mysql□SHOW WARNINGS\G;

***************************1.row***************************

Level:Warning

Code:1048

Message:Column'id'cannot be null

***************************2.row***************************

Level:Warning

Code:1265

Message:Data truncated for column'date'at row 1

2 rows in set(0.00 sec)

mysql□SELECT*FROM a\G;

***************************1.row***************************

id:0

date:0000-00-00

1 row in set(0.00 sec)

---

如果这两个字段被指定为NOT NULL，那么这里的NULL值，以及非法的date类型的数据，例如'2009-02-30'的"日期"对于MySQL来说是不合法的。在这种情况下，只是产生一个warning，对于这些不合法的数据我们应该让它产生一个错误才对。因此MySQL提供了不同的错误级别处理方式，这就需要使用sql_mode，严格地限制某些不合法的数据。

---

mysql□SET sql_mode='STRICT_TRANS_TABLES';

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO a

-□SELECT NULL,'2009-02-30';

ERROR 1048(23000):Column'id'cannot be null

mysql□INSERT INTO a

-□SELECT 1,'2009-02-30';

ERROR 1292(22007):Incorrect date value:'2009-02-30'for column'date'at row 1

---

从上述例子中，sql_mode设置为STRICT_TRANS_TABLES后，对于MySQL来说就是非法数据的插入就会产生错误了。下面一小节将列出比较有用的sql_mode值，对于比较全的介绍读者可以参考MySQL官方手册。

# 4.6.5　ENUM和SET类型

MySQL也支持枚举和集合类型CHECK约束，这就是ENUM和SET类型。很多人认为性别这样的信息，只有male、female这样两种值，适合用ENUM类型进行约束。

---

mysql>CREATE TABLE a(

->id INT,

->sex ENUM('male','female'));

Query OK,0 rows affected(0.12 sec)

mysql>INSERT INTO a

->SELECT 1,'female';

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql>INSERT INTO a

->SELECT 2,'bi';

Query OK,1 row affected,1 warning(0.03 sec)

Records:1 Duplicates:0 Warnings:1

---

## 但从这个例子来看，对于不在枚举范围内的数据，其默认并不是严格的，只有对CHECK约束违反了，才会报错，需要设置sql_mode：

---

mysql>SET sql_mode='STRICT_TRANS_TABLES';

Query OK,0 rows affected(0.00 sec)

mysql>INSERT INTO a

->SELECT 2,'bi';

ERROR 1265(01000):Data truncated for column'sex'at row 1

---

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□CHECK□□□□□□□□□□□□□□□□□□□□□□□□ENUM□SET□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 4.6.6 触发器使用

触发器是一种特殊的存储过程，它在插入、删除或修改特定表中的数据时触发执行，它比数据库本身标准的功能有更精细和更复杂的数据控制能力。

数据库触发器有以下的作用：INSERT、DELETE、UPDATE语句触发它们执行一个SQL查询或者脚本。MySQL 5.0对触发器有了有限的支持，而更加完善的支持会在MySQL 5.1版本中实现。下面将介绍触发器的创建、使用和删除。

创建触发器需要使用CREATE TRIGGER语句，且有Super权限的MySQL用户才可以执行这条语句。

---

CREATE

[DEFINER={user|CURRENT_USER}]

TRIGGER trigger_name BEFORE|AFTER INSERT|UPDATE|DELETE

ON tbl_name FOR EACH ROW trigger_stmt

---

其中触发器名可以有以下6种：操作时间包括INSERT、UPDATE、DELETE、BEFORE、AFTER。操作时间为BEFORE、AFTER，代表触发器是在激活它的语句之前或之后触发。目前，MySQL触发器只支持FOR EACH ROW，即每一行触发，而不像其他的DB2，FOR EACH STATEMENT，即每一句触发。

触发器在某些时候可以增加MySQL的一些新特性，如实现某些形式的引用完整性约束和CHECK约束，不过触发器还是有很多的用途的。下面将介绍触发器的使用方法。

为了演示触发器的使用方法，首先建立一个用户账号表，该表包含"用户编号"和"用户现金"两个字段，要求用户现金的值不能为负数。

---

mysql>CREATE TABLE usercash(

->userid INT NOT NULL,

->cash INT UNSIGNED NOT NULL);

Query OK,0 rows affected(0.11 sec)

mysql›INSERT INTO usercash

-›SELECT 1,1000;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql›UPDATE usercash

-›SET cash=cash-(-20)WHERE userid=1;

Query OK,1 row affected(0.05 sec)

Rows matched:1 Changed:1 Warnings:0

---

# 在本例中，SQL语句看似是在扣除用户账户余额，但实际操作却由于负负得正的原因，反而对用户账户进行了充值操作，这种操作可能会给商家带来损失，下面就利用触发器来实现当发生上述类似操作时进行记录：

---

mysql›CREATE TABLE usercash_err_log(

-›userid INT NOT NULL,

-›old_cash INT UNSIGNED NOT NULL,

-›new_cash INT UNSIGNED NOT NULL,

-›user VARCHAR(30),

-›time DATETIME);

Query OK,0 rows affected(0.13 sec)

mysql›DELIMITER$$

Query OK,0 rows affected(0.00 sec)

mysql›CREATE TRIGGER tgr_usercash_update BEFORE UPDATE ON usercash

-›FOR EACH ROW

-›BEGIN

-›IF new.cash-old.cash›0 THEN

-›INSERT INTO usercash_err_log

-›SELECT old.userid,old.cash,new.cash,USER(),NOW();

-›SET new.cash=old.cash;

-›END IF;

-￼END;

-￼$$

Query OK,0 rows affected(0.00 sec)

mysql￼DELIMITER$$

Query OK,0 rows affected(0.00 sec)

---

为了方便测试，我们先将usercash_err_log表内的记录清空。接下来测试触发器是否工作，触发器tgr_usercash_update是属于在BEFORE触发的，所以每次都会在更新之前判断更新后的现金是否小于零，如果小于，就将要更新的现金值改为历史最小值，并向表cash进行的更改日志记录，即向表usercash_err_log插入日志记录。下面的SQL语句：

---

mysql￼DELETE FROM usercash;

Query OK,1 row affected(0.02 sec)

mysql￼INSERT INTO usercash

-￼SELECT 1,1000;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql￼UPDATE usercash

-￼SET cash=cash-(-20)

-￼WHERE userid=1;

Query OK,0 rows affected(0.02 sec)

Rows matched:1 Changed:0 Warnings:0

mysql￼SELECT*FROM usercash\G;

***************************1.row***************************

userid:1

cash:100

1 row in set(0.00 sec)

mysql￼SELECT*FROM usercash_err_log\G;

***************************1.row***************************

userid:1

old_cash:1000

new_cash:1020

user:root@localhost

time:2009-11-06 11:49:49

Message:Column'id'cannot be null

1 row in set(0.00 sec)

---

可以看到，上面的这个触发器跟踪记录了在向usercash_err_log表中插入数据失败的SQL语句。这里需要提醒读者的是，错误日志信息表中的内容需要人工定期地进行清理，

# 4.6.7 外键约束

外键是否起作用取决于MySQL表的类型，MyISAM类型的表不支持外键，只能起到注释的作用，而对于同时支持事务和外键的InnoDB类型的表而言，则可以设置外键约束，其语法如下：

```
[CONSTRAINT[symbol]]FOREIGN KEY

[index_name](index_col_name,...)

REFERENCES tbl_name(index_col_name,...)

[ON DELETE reference_option]

[ON UPDATE reference_option]

reference_option:

RESTRICT|CASCADE|SET NULL|NO ACTION
```

外键既可以在CREATE TABLE时指定，也可以在表创建以后用ALTER TABLE来指定，下面的例子就说明了如何创建外键：

```
mysql　CREATE TABLE parent(

-　id INT NOT NULL,

-　PRIMARY KEY(id)

-　)ENGINE=INNODB;

Query OK,0 rows affected(0.13 sec)

mysql　CREATE TABLE child(

-　id INT,parent_id INT,

-　FOREIGN KEY(parent_id)REFERENCES parent(id)

-　)ENGINE=INNODB;

Query OK,0 rows affected(0.16 sec)
```

需要注意的是，我们可以在创建外键的同时指定外键的属性，其中ON DELETE、ON UPDATE表示事件触发限制，即在DELETE、UPDATE操作时限制的情况，具体的可设置参数如下所示：

□CASCADE

□SET NULL

□NO ACTION

□RESTRICT

CASCADE□□□□□□DELETE□UPDATE□□□□□□□□□□□□□□□□□□□□□
DELETE□UPDATE□□□SET NULL□□□□□□□□DELETE□UPDATE□□□□□
□□□□□□□□□□□□□□NULL□□□□□□□□□□□□□□□□□□NULL□□NO ACTION□
□□□□□□□DELETE□UPDATE□□□□□□□□□□□□□□□□□□RESTRICT□□□□
□□□□□DELETE□UPDATE□□□□□□□□□□□□□□□□□□□□□□□□□□□□ON
DELETE□ON UPDATE□RESTRICT□□□□□□□□□□□

□□□□□□□□□Oracle□□□□□□□□□□□□□□□deferred check□□□□□□□□□□□
□□SQL□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□immediate
check□□□□□□□□□□□□□□□□□MySQL□□□□NO ACTION□RESTRICT□
□□□□□□□□□

□Oracle□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□
□□□□□□□□□□□□□□□□□□□□□Microsoft SQL Server□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□child□□□□□□□□□□□□□
□□□□□parent_id□□□□□□□□□□□SHOW CREATE TABLE□□□□□
InnoDB□□□□□□□□□□□□□□parent_id□□□□□□□

```
mysql□SHOW CREATE TABLE child\G;

****************************1.row**************************

Table:child

Create Table:CREATE TABLE'child'(

'id'int(11)DEFAULT NULL,

'parent_id'int(11)NOT NULL,

KEY'parent_id'('parent_id'),

CONSTRAINT'child_ibfk_1'FOREIGN KEY('parent_id')REFERENCES'parent'('id')

)ENGINE=InnoDB DEFAULT CHARSET=utf8
```

1 row in set(0.00 sec)

---

对于处于外键约束关系中的表，可以在导入数据之前临时解除这种约束关系，也就是不检查外键约束，这样在MySQL导入数据时就不会考虑先导入哪张表后导入哪张表，从而加快数据导入速度，其具体的做法如下所示。

---

mysql>SET foreign_key_checks=0;

mysql>LOAD DATA……

mysql>SET foreign_key_checks=1;

---

# 4.7 视图

在MySQL中视图又称为View，简单地说视图就是对一条SQL语句结果集的封装，其又被称之为permanent table（虚表），其行为表现和真实的表几乎一样。

## 4.7.1 创建视图

因为视图又被称之为虚表，所以创建视图语句创建出来的可以看成是一张虚拟的表，创建视图所依赖的表被称之为base table（基表）。需要注意的是数据库中只会存储视图的定义，而并不会存储视图所对应的数据。

MySQL在版本5.0之后开始支持视图，其语法格式如下：

---

```
CREATE

[OR REPLACE]

[ALGORITHM={UNDEFINED|MERGE|TEMPTABLE}]

[DEFINER={user|CURRENT_USER}]

[SQL SECURITY{DEFINER|INVOKER}]

VIEW view_name[(column_list)]

AS select_statement

[WITH[CASCADED|LOCAL]CHECK OPTION]
```

---

在上面的语法中只需要重点关注被阴影标注的几个选项即可，其余的选项在使用时出现的频率并不高。其中的算法稍后会详细介绍，可更新视图（updatable view）会在后面介绍，WITH CHECK OPTION也会在后面介绍。现在我们创建一个视图，示例代码如下所示:

---

```
mysql>CREATE TABLE t(id INT);

Query OK,0 rows affected(0.13 sec)

mysql>CREATE VIEW v_t

->AS

->SELECT*FROM t WHERE id<10;
```

Query OK,0 rows affected(0.00 sec)

mysql>INSERT INTO v_t SELECT 20;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql>SELECT*FROM v_t;

Empty set(0.00 sec)

---

上述视图中我们只看到了一个id为10的行在v_t视图中，但是我们插入的是id为20的行，而且这条SQL语句竟然执行成功了。如果我们加上检查选项（WITH CHECK OPTION）呢？

---

mysql>ALTER VIEW v_t

->　AS

->　SELECT*FROM t WHERE id<10

->　WITH CHECK OPTION;

Query OK,0 rows affected(0.00 sec)

mysql>INSERT INTO v_t SELECT 20;

ERROR 1369(HY000):CHECK OPTION failed'mytest.v_t'

---

由于MySQL数据库对于视图的定义并没有所谓的物化视图，因此对视图的查询仅是把视图的定义展开罢了。

MySQL数据库DBA的一个常用的命令是SHOW TABLES，该命令会显示出当前数据库下所有的表。但因为视图是虚表，同样被作为表显示出来，如：

---

mysql>SHOW TABLES\G;

***************************1.row***************************

Tables_in_mytest:t

***************************2.row***************************

Tables_in_mytest：v_t

2 rows in set(0.00 sec)

# 通过SHOW TABLES命令在t表上以v_t查看表也同样属于数据库表。因此，需要人工地在查询命令中筛选出information_schema中存储的TABLE信息，只查询属于BASE TABLE的表。其SQL命令如下：

mysql＞SELECT*FROM information_schema.TABLES

-＞WHERE table_type='BASE TABLE'

-＞AND table_schema=database()\G;

***************************1.row***************************

TABLE_CATALOG:NULL

TABLE_SCHEMA:mytest

TABLE_NAME:t

TABLE_TYPE:BASE TABLE

ENGINE:InnoDB

VERSION:10

ROW_FORMAT:Compact

TABLE_ROWS:1

AVG_ROW_LENGTH:16384

DATA_LENGTH:16384

MAX_DATA_LENGTH:0

INDEX_LENGTH:0

DATA_FREE:0

AUTO_INCREMENT:NULL

CREATE_TIME:2009-11-09 16:27:52

UPDATE_TIME:NULL

CHECK_TIME:NULL

TABLE_COLLATION:utf8_general_ci

CHECKSUM:NULL

CREATE_OPTIONS:

TABLE_COMMENT:

1 row in set(0.00 sec)

视图相对于表而言有很多的meta data信息，通过information_schema库的VIEWS可以查看视图的详细定义信息，包括视图的definer信息、安全属性、字符集等信息，下面来看一下VIEWS的属性。

mysql>SELECT*FROM

->information_schema.VIEWS

->WHERE table_schema=database()\G;

***************************1.row***************************

TABLE_CATALOG:NULL

TABLE_SCHEMA:mytest

TABLE_NAME:v_t

VIEW_DEFINITION:select'mytest'.'t'.'id'AS'id'from'mytest'.'t'where('mytest'.'t'.'id'<10)

CHECK_OPTION:CASCADED

IS_UPDATABLE:YES

DEFINER:root@localhost

SECURITY_TYPE:DEFINER

CHARACTER_SET_CLIENT:latin1

COLLATION_CONNECTION:latin1_swedish_ci

1 row in set(0.00 sec)

## 4.7.2 刷新模式

Oracle具体化视图的定义——我们将在后面详细说明——可以包括将查询结果存储在数据库中所需的所有查询功能。具体而言，它们可以包括任意复杂的查询以及JOIN连接或者GROUP BY聚合。通过这些SQL特性，可以将许多表连接到一起，或者将数据聚合成汇总信息并把结果保存在具体化视图中。正因为如此，具体化视图的刷新比Microsoft SQL Server中的索引视图要复杂得多。

在Oracle中，在创建具体化视图时有两种生成选项：

❑BUILD IMMEDIATE

❑BUILD DEFERRED

BUILD IMMEDIATE选项表示在创建具体化视图时立即将数据写入视图，而BUILD DEFERRED选项则将数据的生成推迟到某个后续时间再进行。

具体化视图中的数据从定义该视图的基表中获得，该视图中的数据需要偶尔刷新，以便与基表中可能发生的变更保持同步。虽然该视图的SQL定义在创建视图时就已经固定，但是视图中的数据将随时间变化。

视图的刷新是与源表发生的DML操作分开进行的，但是可以按几种不同的方式予以控制：

❑ON DEMAND

❑ON COMMIT

ON DEMAND表示具体化视图将按某个计划或请求进行刷新，ON COMMIT表示每当基表中提交了对数据的DML操作时，该视图都将进行刷新。

可以使用如下刷新方法：

❑FAST

❑COMPLETE

❑FORCE

□NEVER

FAST□□□□□□□□□□□□□□□□□□□□□□□□□□□□COMPLETE□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□FORCE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
FAST□□□□□□COMPLETE□□□□NEVER□□□□□□□□□□□□□□□

MySQL□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□ON DEMAND□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□CREATE TABLE Orders

-□(

-□order_id INT UNSIGNED NOT NULL AUTO_INCREMENT,

-□product_name VARCHAR(30)NOT NULL,

-□price DECIMAL(8,2)NOT NULL,

-□amount SMALLINT NOT NULL,

-□PRIMARY KEY(order_id)

-□)ENGINE=InnoDB;

Query OK,0 rows affected(0.13 sec)

mysql□INSERT INTO Orders VALUES

-□(NULL,'CPU',135.5,1),

-□(NULL,'Memory',48.2,3),

-□(NULL,'CPU',125.6,3),

-□(NULL,'CPU',105.3,4)

-□;

Query OK,4 rows affected(0.03 sec)

Records:4 Duplicates:0 Warnings:0

mysql□SELECT*FROM Orders\G;

***************************1.row*************************

order_id:1

product_name:CPU

price:135.50

amount:1

***************************2.row***************************

order_id:2

product_name:Memory

price:48.20

amount:3

***************************3.row***************************

order_id:3

product_name:CPU

price:125.60

amount:3

***************************4.row***************************

order_id:4

product_name:CPU

price:105.30

amount:4

4 rows in set(0.00 sec)

---

# 接下来创建视图表。为了演示后面的技巧，这里创建的是真实表

---

mysql>CREATE TABLE Orders_MV(

->product_name VARCHAR(30)NOT NULL

->,price_sum DECIMAL(8,2)NOT NULL

->,amount_sum INT NOT NULL

->,price_avg FLOAT NOT NULL

->,orders_cnt INT NOT NULL

->,UNIQUE INDEX(product_name)

->);

Query OK,0 rows affected(0.13 sec)

mysql□INSERT INTO Orders_MV

-□SELECT product_name

-□,SUM(price),SUM(amount),AVG(price)

-□,COUNT(*)

-□FROM Orders

-□GROUP BY product_name;

Query OK,2 rows affected(0.02 sec)

Records:2 Duplicates:0 Warnings:0

mysql□

mysql□

mysql□SELECT*FROM Orders_MV\G;

***************************1.row***************************

product_name:CPU

price:366.40

amount_sum:8

price_avg:122.133

orders_cnt:3

***************************2.row***************************

product_name:Memory

price:48.20

amount_sum:3

price_avg:48.2

orders_cnt:1

2 rows in set(0.00 sec)

---

□□□□□□□□□□□□□□□□□□□□□□□□□Orders_MV□□□□_MV□□□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ON DEMAND□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□COMPLETE□□□□□□□□□□□□FAST□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□order_id□□□□□

由于无须支持增量式ON COMMIT类型的物化视图的创建，因此可以将Oracle数据库系统中的物化视图创建语句转换到MySQL数据库系统中的触发器实现，对应的增量式物化视图触发器，针对Orders表的插入操作创建如下：

```
DELIMITER$$

CREATE TRIGGER tgr_Orders_insert

AFTER INSERT ON Orders

FOR EACH ROW

BEGIN

SET@old_price_sum=0;

SET@old_amount_sum=0;

SET@old_price_avg=0;

SET@old_orders_cnt=0;

SELECT IFNULL(price_sum,0),IFNULL(amount_sum,0),IFNULL(price_avg,0),IFNULL(orders_cnt,0)

FROM Orders_MV

WHERE product_name=NEW.product_name

INTO@old_price_sum,@old_amount_sum,@old_price_avg,@old_orders_cnt;

SET@new_price_sum=@old_price_sum+NEW.price;

SET@new_amount_sum=@old_amount_sum+NEW.amount;

SET@new_orders_cnt=@old_orders_cnt+1;

SET@new_price_avg=@new_price_sum/@new_orders_cnt;

REPLACE INTO Orders_MV

VALUES(NEW.product_name,@new_price_sum,@new_amount_sum,@new_price_avg,@new_orders_cnt);

END;

$$

DELIMITER;
```

在上面的触发器中，当INSERT操作发生时，先根据INSERT操作的数据查询一次Orders_MV物化视图表，然后对查询的结果进行计算，最后将计算的结果更新到Orders_MV物化视图表

```
mysql>INSERT INTO Orders VALUES(NULL,'SSD',299,3);

Query OK,1 row affected,1 warning(0.03 sec)

mysql>INSERT INTO Orders VALUES(NULL,'Memory',47.9,5);

Query OK,1 row affected(0.03 sec)

mysql>SELECT*FROM Orders_MV\G;

***************************1.row***************************

product_name:CPU

price:366.40

amount_sum:8

price_avg:122.133

orders_cnt:3

***************************2.row***************************

product_name:Memory

price:96.10

amount_sum:8

price_avg:48.05

orders_cnt:2

***************************3.row***************************

product_name:SSD

price:299.00

amount_sum:3

price_avg:299

orders_cnt:1

3 rows in set(0.00 sec)
```

从上述执行结果中可以看出，对于物化视图Orders_MV，在执行查询操作时并没有涉及聚合函数等复杂的SQL语句。这主要是因为在使用ON_COMMIT方式定义物化视图时，每当对基表Orders执行插入、UPDATE、DELETE操作时，或者对基表执行DELETE、UPDATE操作后，物化视图便会自动更新。

□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Query Rewrite□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 4.8  分区表

## 4.8.1  分区概述

分区功能并不是在存储引擎层完成的，因此不是只有InnoDB存储引擎支持分区，常见的存储引擎MyISAM、NDB等都支持。但也并不是所有的存储引擎都支持，如CSV、FEDORATED、MERGE等就不支持分区。在使用分区功能前，应该对自己使用的存储引擎对分区的支持有所了解。

MySQL数据库在5.1版本时添加了对分区的支持。分区的过程是将一个表或索引分解为多个更小、更可管理的部分。就访问数据库的应用而言，从逻辑上讲，只有一个表或一个索引，但是在物理上这个表或索引可能由数十个物理分区组成。每个分区都是独立的对象，可以独自处理，也可以作为一个更大对象的一部分进行处理。

MySQL数据库支持的分区类型为水平分区 [1]，并不支持垂直分区 [2]。此外，MySQL数据库的分区是局部分区索引，一个分区中既存放了数据又存放了索引。而全局分区是指，数据存放在各个分区中，但是所有数据的索引放在一个对象中。目前，MySQL数据库还不支持全局分区。

是否启用了分区功能可以通过如下命令来观察：

---

```
mysql＞SHOW VARIABLES LIKE'%partition%'\G;
```

```
***************************1.row***************************
```

```
Variable_name:have_partitioning
```

```
Value:YES
```

```
1 row in set(0.00 sec)
```

---

分区功能也可以通过SHOW PLUGINS来查看：

---

```
mysql＞SHOW PLUGINS\G;
```

```
……
```

```
***************************2.row***************************
```

```
Name:partition
```

```
Status:ACTIVE
```

```
Type:STORAGE ENGINE
```

Library:NULL

License:GPL

……

9 rows in set(0.01 sec)

---

□□□DBA□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLTP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□MySQL□□□□□□□□□□□□□□□□□□□□

□RANGE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL 5.5□□□□□RANGE COLUMNS□□□□□

□LIST□□□□□RANGE□□□□□□□□□□LIST□□□□□□□□□□□□□□MySQL 5.5□□□□□LIST COLUMNS□□□□□

□HASH□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□KEY□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□

---

mysql□CREATE TABLE t1(

-□col1 INT NOT NULL,

-□col2 DATE NOT NULL,

-□col3 INT NOT NULL,

-□col4 INT NOT NULL,

-□UNIQUE KEY(col1,col2)

-□)

-□PARTITION BY HASH(col3)

-□PARTITIONS 4;

ERROR 1503(HY000):A PRIMARY KEY must include all columns in the table's partitioning function

## 如果我们将主键改为NULL，取消主键约束，这个时候就可以了，因为不再有主键约束的限制，这样也就可以实现了

mysql〉CREATE TABLE t1(

-〉col1 INT NULL,

-〉col2 DATE NULL,

-〉col3 INT NULL,

-〉col4 INT NULL,

-〉UNIQUE KEY(col1,col2,col3,col4)

-〉)

-〉PARTITION BY HASH(col3)

-〉PARTITIONS 4;

Query OK,0 rows affected(0.53 sec)

## 同样对于有多个唯一索引，但是又不符合唯一索引的要求下，我们还是要能进行分区的，那么我们的SQL语句我们还是要这样子写

CREATE TABLE t1(

col1 INT NULL,

col2 DATE NULL,

col3 INT NULL,

col4 INT NULL

)engine=innodb

PARTITION BY HASH(col3)

PARTITIONS 4;

CREATE TABLE t1(

col1 INT NULL,

col2 DATE NULL,

```
col3 INT NULL,

col4 INT NULL,

key(col4)

)engine=innodb

PARTITION BY HASH(col3)

PARTITIONS 4;
```

---

[1] 如果表中存在主键或唯一键时，分区键必须是包含主键或唯一键的列

[2] 如果没有主键或唯一键，分区键可以是任意列，但此时还有必须是整数类型

# 4.8.2 □□□□

## 1.RANGE□□

□□□□□□□□□□□□□□□RANGE□□□□□□□□□□□□□□□□□□□□□□CREATE TABLE □□□□□□□□□□id□□□□□□□□□□id□□10□□□□□□□p0□□□□□id□□□□10□□20□□□□ □□□□p1□□□□

---

CREATE TABLE t(

id INT

)ENGINE=INNDB

PARTITION BY RANGE(id)(

PARTITION p0 VALUES LESS THAN(10),

PARTITION p1 VALUES LESS THAN(20));

---

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ibd□□□□□□□□□□□□□□□□□□□□□□□□ibd □□□□□□□□□□t#P#p0.ibd□t#P#p1.ibd□

---

mysql□system ls-lh/usr/local/mysql/data/test2/t*

-rw-rw----1 mysql mysql 8.4K 7□31 14:11/usr/local/mysql/data/test2/t.frm

-rw-rw----1 mysql mysql 28 7□31 14:11/usr/local/mysql/data/test2/t.par

-rw-rw----1 mysql mysql 96K 7□31 14:12/usr/local/mysql/data/test2/t#P#p0.ibd

-rw-rw----1 mysql mysql 96K 7□31 14:12/usr/local/mysql/data/test2/t#P#p1.ibd

---

□□□□□□□□□□□□□□

---

mysql□INSERT INTO t SELECT 9;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO tSELECT 10;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql：INSERT INTO t SELECT 15;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

---

## 对于表t，由于id列插入了10、15这两个值，我们需要查看表在不同分区中数据的分布情况，可以通过查询information\_schema架构下的PARTITIONS表来进行，如下面语句所示：

---

mysql：SELECT*FROM information_schema.PARTITIONS

-　WHERE table_schema=database()AND table_name='t'\G;

***************************1.row***************************

TABLE_CATALOG:NULL

TABLE_SCHEMA:test2

TABLE_NAME:t

PARTITION_NAME:p0

SUBPARTITION_NAME:NULL

PARTITION_ORDINAL_POSITION:1

SUBPARTITION_ORDINAL_POSITION:NULL

PARTITION_METHOD:RANGE

SUBPARTITION_METHOD:NULL

PARTITION_EXPRESSION:id

SUBPARTITION_EXPRESSION:NULL

PARTITION_DESCRIPTION:10

TABLE_ROWS:1

AVG_ROW_LENGTH:16384

DATA_LENGTH:16384

MAX_DATA_LENGTH:NULL

INDEX_LENGTH:0

DATA_FREE:0

CREATE_TIME:NULL

UPDATE_TIME:NULL

CHECK_TIME:NULL

CHECKSUM:NULL

PARTITION_COMMENT:

NODEGROUP:default

TABLESPACE_NAME:NULL

*************************2.row*************************

TABLE_CATALOG:NULL

TABLE_SCHEMA:test2

TABLE_NAME:t

PARTITION_NAME:p1

SUBPARTITION_NAME:NULL

PARTITION_ORDINAL_POSITION:2

SUBPARTITION_ORDINAL_POSITION:NULL

PARTITION_METHOD:RANGE

SUBPARTITION_METHOD:NULL

PARTITION_EXPRESSION:id

SUBPARTITION_EXPRESSION:NULL

PARTITION_DESCRIPTION:20

TABLE_ROWS:2

AVG_ROW_LENGTH:8192

DATA_LENGTH:16384

MAX_DATA_LENGTH:NULL

INDEX_LENGTH:0

DATA_FREE:0

CREATE_TIME:NULL

UPDATE_TIME:NULL

CHECK_TIME:NULL

CHECKSUM:NULL

PARTITION_COMMENT:

NODEGROUP:default

TABLESPACE_NAME:NULL

2 rows in set(0.00 sec)

---

TABLE_ROWS列反映了每个分区中记录的数量。由于测试表中插入了9、10、15三条记录，因此分区p0中有1条记录、分区p1中有2条记录。PARTITION_METHOD表示分区的类型，这里显示为RANGE。

对于表t，由于我们定义的分区是有限的，如果插入分区范围外的数据就会报错，MySQL数据库会抛出异常。如这里向表t中插入30这条记录：

---

mysql□INSERT INTO t SELECT 30;

ERROR 1526(HY000):Table has no partition for value 30

---

对于上述问题，我们可以对分区添加一个MAXVALUE值的分区。MAXVALUE可以理解为正无穷，因此所有大于等于20且小于MAXVALUE的值放入p2分区。

---

mysql□ALTER TABLE t

-□ADD PARTITION(

-□partition p2 values less than maxvalue);

Query OK,0 rows affected(0.45 sec)

Records:0 Duplicates:0 Warnings:0

mysql□INSERT INTO t SELECT 30;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

---

RANGE分区主要用于日期列的分区，例如对于销售类的表，可以根据销售的时间进行分区，如表sales：

---

mysql□CREATE TABLE sales(

```
-　money INT UNSIGNED NOT NULL,

-　date DATETIME

-　)ENGINE=INNODB

-　PARTITION by RANGE(YEAR(date))(

-　PARTITION p2008 VALUE LESS THEN(2009),

-　PARTITION p2009 VALUE LESS THEN(2010),

-　PARTITION p2010 VALUE LESS THEN(2011)

-　);

Query OK,0 rows affected(0.34 sec)

mysql　INSERT INTO sales SELECT 100,'2008-01-01';

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql　INSERT INTO sales SELECT 100,'2008-02-01';

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql　INSERT INTO sales SELECT 200,'2008-01-02';

Query OK,1 row affected(0.04 sec)

Records:1 Duplicates:0 Warnings:0

mysql　INSERT INTO sales SELECT 100,'2009-03-01';

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql　INSERT INTO sales SELECT 200,'2010-03-01';

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0
```

此时，如果我们要删除sales表中所有关于年份在2008年的数据，则我们可以用DELETE FROM sales WHERE date　='2008-01-01'and date　'2009-01-01'，删除所有2008年的数据，也可以用：

```
mysql　alter table sales drop partition p2008;
```

Query OK,0 rows affected(0.18 sec)

Records:0 Duplicates:0 Warnings:0

---

在查询数据时可以通过查看分区的信息来判断数据是否落在2008年区间。具体的操作步骤如下：

---

mysql□EXPLAIN PARTITIONS

-□SELECT*FROM sales

-□WHERE date□='2008-01-01'AND date□='2008-12-31'\G;

***************************1.row***************************

id:1

select_type:SIMPLE

table:sales

partitions:p2008

type:ALL

possible_keys:NULL

key:NULL

key_len:NULL

ref:NULL

rows:5

Extra:Using where

1 row in set(0.00 sec)

---

通过EXPLAIN PARTITION语句可以查看到，指定查询范围的SQL语句只落在分区p2008，这也是分区表带来的好处之一——称为Partition Pruning。所谓分区修剪就是在查询中数据库优化器只扫描那些特定分区，而不是扫描所有的分区，这样就可以在很大程度上提高查询的效率。

---

mysql□EXPLAIN PARTITIOENS

-□SELECT*FROM sales

-□WHERE date□='2008-01-01'AND date□'2009-01-01'\G;

```
*************************1.row*************************

id:1

select_type:SIMPLE

table:sales

partitions:p2008,p2009

type:ALL

possible_keys:NULL

key:NULL

key_len:NULL

ref:NULL

rows:5

Extra:Using where

1 row in set(0.00 sec)
```

---

# 很明显当满足date＜'2009-01-01'并且date＞='2008-12-31'的时候，数据就只存在于p2008和p2009两个分区中了，这样就避免了对其他分区数据的扫描，提高了SQL的效率。

## 创建sales表，要求按照年月进行分区，例如：本表中有三个分区，分别存储三个月的数据。

---

```
mysql□CREATE TABLE sales(

-□money INT UNSIGNED NOT NULL,

-□date DATETIME

-□)ENGINE=INNODB

-□PARTITION by RANGE(YEAR(date)*100+MONTH(date))(

-□PARTITION p201001 VALUES LESS THEN(201002),

-□PARTITION p201002 VALUES LESS THEN(201003),

-□PARTITION p201003 VALUES LESS THEN(201004)

-□);

Query OK,0 rows affected(0.37 sec)
```

---

根据以上的SQL语句看到分区后每个分区代表的范围。我们再执行下面的SQL语句来看一下分区的执行情况：

```
mysql＞EXPLAIN PARTITIONS

-＞SELECT*FROM sales

-＞WHERE date＞='2010-01-01'AND date＜='2010-01-31'\G;

***************************1.row***************************

id:1

select_type:SIMPLE

table:sales

partitions:p201001,p201002,p201003

type:ALL

possible_keys:NULL

key:NULL

key_len:NULL

ref:NULL

rows:4

Extra:Using where

1 row in set(0.00 sec)
```

可以看到它分别扫描了p201001、p201002、p201003三个分区的数据，然后进行合并。在创建RANGE分区时，我们可以使用函数YEAR()、TO_DAYS()、TO_SECONDS()、UNIX_TIMESTAMP()，来将日期转化为整数值，然后进行分区。下面给出了使用TO_DAYS的例子：

```
mysql＞CREATE TABLE sales(

-＞money INT UNSIGNED NOT NULL,

-＞date DATETO,E

-＞)ENGINE=INNODB

-＞PARTITION by range(TO_DAYS(date))(
```

-　PARTITION p201001

-　VALUES LESS THEN(TO_DAYS('2010-02-01')),

-　PARTITION p201002

-　VALUES LESS THEN(TO_DAYS('2010-03-01')),

-　PARTITION p201003

-　VALUES LESS THEN(TO_DAYS('2010-04-01'))

-　);

Query OK,0 rows affected(0.36 sec)

## 查询数据时，可以通过以下方式，查看数据处于哪个分区当中：

mysql　EXPLAIN PATITIONS

-　SELECT*FROM sales

-　WHERE date　='2010-01-01'AND date　='2010-01-31'\G;

***************************1.row***************************

id:1

select_type:SIMPLE

table:sales

partitions:p201001

type:ALL

possible_keys:NULL

key:NULL

key_len:NULL

ref:NULL

rows:4

Extra:Using where

1 row in set(0.00 sec)

## 2.LIST分区

# LIST□□□RANGE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
mysql□CREATE TABLE t(

-□a INT,

-□b INT)ENGINE=INNODB

-□PARTITION BY LIST(b)(

-□PARTITION p0 VALUES IN(1,3,5,7,9),

-□PARTITION p1 VALUES IN(0,2,4,6,8)

-□);

Query OK,0 rows affected(0.26 sec)
```

## □□□RANGE□□□□□□□□VALUES LESS THAN□□□LIST□□□□VALUES IN□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□t□□□□□□□□□□

```
mysql□INSERT INTO t SELECT 1,1;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO t SELECT 1,2;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO t SELECT 1,3;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO t SELECT 1,4;

Query OK,1 row affected(0.03 sec)

Records:1 Duplicates:0 Warnings:0

mysql□SELECT table_name,partition_name,table_rows

-□FROM information_schema.PARTITIONS

-□WHERE table_name='t'AND table_schema=DATABASE()\G;

***************************1.row***************************
```

table_name:t

partition_name:p0

table_rows:2

***************************2.row***************************

table_name:t

partition_name:p1

table_rows:2

2 rows in set(0.00 sec)

---

# 对于其他分区类型的数据库，MySQL数据库同样也是通过

---

mysql〉INSERT INTO t SELECT 1,10;

ERROR 1526(HY000):Table has no partition for value 10

---

# 在使用INSERT插入多个行数据的过程中遇到分区未定义的值时，MyISAM和InnoDB存储引擎的处理完全不同。MyISAM引擎会将之前的行数据都插入，但之后的数据不会被插入。而InnoDB存储引擎将其视为一个事务，因此没有任何数据被插入。下面来显示这两者的区别。MyISAM存储引擎下的操作

---

mysql〉CRATE TABLE t(

-〉a INT,

-〉b INT)ENGINE=MyISAM

-〉PARTITION BY LIST(b)(

-〉PARTITION p0 VALUES IN(1,3,5,7,9),

-〉PARTITION p1 VALUES IN(0,2,4,6,8)

-〉);

Query OK,0 rows affected(0.05 sec)

mysql〉INSERT INTO t VALUES(1,2),(2,4),(6,10),(5,3);

ERROR 1526(HY000):Table has no partition for value 10

mysql〉SELECT*FROM t;

+------+------+

|a|b|

+------+------+

|1|2|

|2|4|

+------+------+

2 rows in set(0.00 sec)

---

从上面可以看出6、10这行和5、3这行数据被筛选掉了，1、2这行和2、4这行数据被正确保存下来了。如果要清除数据，可以使用InnoDB引擎重建分区，如下：

---

mysql＞TRUNCATE TABLE t;

Query OK,2 rows affected(0.00 sec)

mysql＞ALTER TABLE t ENGINE=InnoDB;

Query OK,0 rows affected(0.25 sec)

Records:0 Duplicates:0 Warnings:0

mysql＞INSERT INTO t VALUES(1,2),(2,4),(6,10),(5,3);

ERROR 1526(HY000):Table has no partition for value 10

mysql＞SELECT*FROM t;

Empty set(0.00 sec)

---

从上面可以看出，因为6、10这行记录无法找到匹配的分区，使得表t的所有数据都没有被保存下来，所以上面查询数据时返回的是空值。

## 3.HASH分区

HASH分区主要用来确保数据在预先确定数目的分区中平均分布，在分区中，只需要基于将要进行哈希的列值指定一个列值或表达式，以及指定被分区的表将要被分割成的分区数量。使用HASH分区时，MySQL会根据用户自定义的表达式的返回值来直接分配分区的数量，而这个返回值不能为负数。

要使用HASH分区来分割一个表，要在CREATE TABLE语句上添加一个"PARTITION BY HASH（expr）"子句，其中"expr"是一个返回一个整数的表达式。它可以仅仅是字段类型

MySQL数据库还支持在创建表时用“PARTITIONS num”语句来告诉数据库创建num个分区，每个分区可能保存的数据量由应用决定，且PARTITIONS关键字紧跟在分区方法之后。如1：

以下创建了一个根据HASH分区的表t，分区列为列b，如下：

---

CREATE TABLE t_hash(

a INT,

b DATETIME

)ENGINE=InnoDB

PARTITION BY HASH(YEAR(b))

PARTITIONS 4;

---

如果插入一个列b为2010-04-01的记录到表t_hash中，那么保存该条记录的分区如下：

---

MOD(YEAR('2010-04-01')，4)

=MOD(2010,4)

=2

---

即记录会放入分区p2中，验证如下所示，且记录放入如下所示：

---

mysql＞INSERT INTO t_hash SELECT 1,'2010-04-01';

Query OK,1 row affected(0.04 sec)

Records:1 Duplicates:0 Warnings:0

mysql＞SELECT table_name,partition_name,table_rows

-＞FROM information_schema.PARTITIONS

-＞WHERE table_schema=DATABASE()AND table_name='t_hash'\G;

***************************1.row***************************

table_name:t_hash

partition_name:p0

table_rows:0

***************************2.row***************************

table_name:t_hash

partition_name:p1

table_rows:0

***************************3.row***************************

table_name:t_hash

partition_name:p2

table_rows:1

***************************4.row***************************

table_name:t_hash

partition_name:p3

table_rows:0

4 rows in set(0.00 sec)

---

可以看到p2包含了1条记录，读者需要注意的是，当没有与分区键匹配的值时，记录将会根据YEAR函数，计算出实际应该对应的分区，然后存放到正确的HASH分区中。这也是为什么分区的意义，即根据分区键将记录分散存放的意义。

MySQL数据库还支持一种LINEAR HASH的分区，它使用一个更加复杂的算法来确定新行插入到已经分区的表中的位置。它的语法和HASH分区的语法相识，只是将关键字HASH改为LINEAR HASH。下面创建一个LINEAR HASH分区的表t_linear_hash，它和表t_hash的区别仅仅在于分区类型的不同。

---

```
CREATE TABLE t_linear_hash(

a INT,

b DATETIME

)ENGINE=InnoDB

PARTITION BY LINEAR HASH(YEAR(b))

PARTITIONS 4;
```

---

如插入行'2010-04-01'，首先由于MySQL将分区的数量向上取为最接近的二次幂值，4的下一个2次幂是V，V=POWER(2，CEILING(LOG(2，num)))=4。

然后计算N=YEAR('2010-04-01')&(V-1)=2。

因此记录保存到P2分区，这就是上面看到的线性HASH分区表，在插入数据后数据变更的详细情况。

---

mysql＞INSERT INTO t_linear_hash SELECT 1,'2010-04-01';

Query OK,1 row affected(0.02 sec)

Records:1 Duplicates:0 Warnings:0

mysql＞SELECT table_name,partition_name,table_rows

-＞FROM information_schema.PARTITIONS

-＞WHERE table_schema=DATABASE()

-＞AND table_name='t_linear_hash'\G;

***************************1.row**************************

table_name:t_linear_hash

partition_name:p0

table_rows:0

***************************2.row**************************

table_name:t_linear_hash

partition_name:p1

table_rows:0

***************************3.row**************************

table_name:t_linear_hash

partition_name:p2

table_rows:1

***************************4.row**************************

table_name:t_linear_hash

partition_name:p3

table_rows:0

4 rows in set(0.01 sec)

---

LINEAR HASH□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□HASH□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 4.KEY□□

KEY□□□□HASH□□□□□□□□□□□□□□HASH□□□□□□□□□□□□□□□□□□□□□□□KEY□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□NDB Cluster□□□□MySQL□□□□□□□MD5□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□PASSWORD()□□□□□□□□□□□□□

---

```
mysql□CREATE TABLE t_key(
-□a INT,
-□b DATETIME)ENGINE=InnoDB
-□PARTITION BY KEY(b)
-□PARTITIONS 4;
Query OK,0 rows affected(0.43 sec)
```

---

□KEY□□□□□□□□□□□□LINEAR□□HASH□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□2□□□□powers-of-two□□□□□□□□□□□□□□□□□□□□□□□□□

## 5.COLUMNS□□

□□□□□□□□RANGE□LIST□HASH□KEY□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□interger□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□YEAR()□TO_DAYS()□MONTH()□□□□□MySQL5.5□□□□□□□□COLUMNS□□□□□□□□□RANGE□□□□LIST□□□□□□□□□□COLUMNS□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□RANGECOLUMNS□□□□□□□□□□□□□□□□□□□

COLUMNS□□□□□□□□□□□□□□□□□

☐所有整数类型，例如INT、SMALLINT、TINYINT、BIGINT、FLOAT和DECIMAL等等。

☐日期类型，例如DATE、DATETIME。不支持其他的日期类型。

☐字符串类型，例如CHAR、VARCHAR、BINARY、VARBINARY、BLOB、TEXT等等这些类型。

与第一种类型相比，基于列的分区不需要使用YEAR()、TO_DAYS()等函数进一步的转换。基于列的COLUMNS的方式。

CREATE TABLE t_columns_range(

a INT,

b DATETIME

)ENGINE=INNODB

PARTITION BY RANGE COLUMNS(B)(

PARTITION p0 VALUES LESS THAN('2009-01-01'),

PARTITION p1 VALUES LESS THAN('2010-01-01')

);

## 多个列的分区也是非常简单的：

CREATE TABLE customers_1(

first_name VARCHAR(25),

last_name VARCHAR(25),

street_1 VARCHAR(30),

street_2 VARCHAR(30),

city VARCHAR(15),

renewal DATE

)

PARTITION BY LIST COLUMNS(city)(

```
PARTITION pRegion_1

VALUES IN('Oskarshamn','Högsby','Mönsterås'),

PARTITION pRegion_2

VALUES IN('Vimmerby','Hultsfred','Västervik'),

PARTITION pRegion_3

VALUES IN('Nässjö','Eksjö','Vetlanda'),

PARTITION pRegion_4

VALUES IN('Uppvidinge','Alvesta','Växjo')

);
```

## 多列RANGE COLUMNS分区，即包含多个列作为分区键：

```
CREATE TABLE rcx(

a INT,

b INT,

c CHAR(3),

d INT

)Engine=InnoDB

PARTITION BY RANGE COLUMNS(a,d,c)(

PARTITION p0 VALUES LESS THAN(5,10,'ggg'),

PARTITION p1 VALUES LESS THAN(10,20,'mmmm'),

PARTITION p2 VALUES LESS THAN(15,30,'sss'),

PARTITION p3 VALUES LESS THAN(MAXVALUE,MAXVALUE,MAXVALUE)

);
```

MySQL5.5开始支持COLUMNS分区，是对原来RANGE、LIST分区的一种发展，RANGE COLUMNS、LIST COLUMNS支持非整形的分区。

# 4.8.3 　子分区

子分区（subpartitioning）是在分区的基础上再进行分区，有时也称这种分区为复合分区（composite partitioning）。MySQL数据库允许在RANGE和LIST的分区上再进行HASH或KEY的子分区，如：

---

```
mysql＞CREATE TABLE ts(a INT,b DATE)engine=innodb

-＞PARTITION BY RANGE(YEAR(b))

-＞SUBPARTITION BY HASH(TO_DAYS(b))

-＞SUBPARTITIONS 2(

-＞PARTITION p0 VALUES LESS THAN(1990),

-＞PARTITION p1 VALUES LESS THAN(2000),

-＞PARTITION p2 VALUES LESS THAN MAXVALUE

-＞);

Query OK,0 rows affected(0.01 sec)

mysql＞system ls-lh/usr/local/mysql/data/test2/ts*

-rw-rw----1 mysql mysql 8.4K Aug 1 15:50/usr/local/mysql/data/test2/ts.frm

-rw-rw----1 mysql mysql 96 Aug 1 15:50/usr/local/mysql/data/test2/ts.par

-rw-rw----1 mysql mysql 96K Aug 1 15:50/usr/local/mysql/data/test2/ts#P#p0#SP#p0sp0.ibd

-rw-rw----1 mysql mysql 96K Aug 1 15:50/usr/local/mysql/data/test2/ts#P#p0#SP#p0sp1.ibd

-rw-rw----1 mysql mysql 96K Aug 1 15:50/usr/local/mysql/data/test2/ts#P#p1#SP#p1sp0.ibd

-rw-rw----1 mysql mysql 96K Aug 1 15:50/usr/local/mysql/data/test2/

ts#P#p1#SP#p1sp1.ibd

-rw-rw----1 mysql mysql 96K Aug 1 15:50/usr/local/mysql/data/test2/ts#P#p2#SP#p2sp0.ibd

-rw-rw----1 mysql mysql 96K Aug 1 15:50/usr/local/mysql/data/test2/ts#P#p2#SP#p2sp1.ibd
```

---

表ts先根据b列进行了RANGE分区，然后又进行了一次HASH分区，所以分区的数量应该为（3×2=）6个，这通过上面的文件可以得到证实。当然我们也可以通过在每个SUBPARTITION语法上来显式地指出各个子分区的名字，如对于表ts，可以这样创建：

---

```
mysql⬚CREATE TABLE ts(a INT,b DATE)

-⬚PARTITION BY RANGE(YEAR(b))

-⬚SUBPARTITION BY HASH(TO_DAYS(b))(

-⬚PARTITION p0 VALUES LESS THAN(1990)(

-⬚SUBPARTITION s0,

-⬚SUBPARTITION s1

-⬚),

-⬚PARTITION p1 VALUES LESS THAN(2000)(

-⬚SUBPARTITION s2,

-⬚SUBPARTITION s3

-⬚),

-⬚PARTITION p2 VALUES LESS THAN MAXVALUE(

-⬚SUBPARTITION s4,

-⬚SUBPARTITION s5

-⬚)

-⬚);

Query OK,0 rows affected(0.00 sec)
```

---

⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚

⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚

⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚SUBPARTITION⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚

---

```
mysql⬚CREATE TABLE ts(a INT,b DATE)

-⬚PARTITION BY RANGE(YEAR(b))

-⬚SUBPARTITION BY HASH(TO_DAYS(b))(

-⬚PARTITION p0 VALUES LESS THAN(1990)(

-⬚SUBPARTITION s0,

-⬚SUBPARTITION s1
```

-□),

-□PARTITION p1 VALUES LESS THAN(2000),

-□PARTITION p2 VALUES LESS THAN MAXVALUE(

-□SUBPARTITION s2,

-□SUBPARTITION s3

-□)

-□);

ERROR 1064(42000):Wrong number of subpartitions defined,mismatch with previous setting near'

PARTITION p2 VALUES LESS THAN MAXVALUE(

SUBPARTITION s2,

SUBPARTITION s3

)

)'at line 8

---

# □□□SUBPARTITION□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□CREATE TABLE ts(a INT,b DATE)

-□PARTITION BY RANGE(YEAR(b))

-□SUBPARTITION BY HASH(TO_DAYS(b))(

-□PARTITION p0 VALUES LESS THAN(1990)(

-□SUBPARTITION s0,

-□SUBPARTITION s1

-□),

-□PARTITION p1 VALUES LESS THAN(2000)(

-□SUBPARTITION s0,

-□SUBPARTITION s1

-□),

-□PARTITION p2 VALUES LESS THAN MAXVALUE(

-□SUBPARTITION s0,

-□SUBPARTITION s1

-□)

-□);

ERROR 1517(HY000):Duplicate partition name s0

---

# 每个分区和子分区都可以指定一个数据目录和索引目录。下面的例子使用6个磁盘，分别是/disk0、/disk1、/disk2……具体的建表语句如下所示：

---

mysql□CREATE TABLE ts(a INT,b DATE)ENGINE=MYISAM

-□PARTITION BY RANGE(YEAR(b))

-□SUBPARTITION BY HASH(TO_DAYS(b))(

-□PARTITION p0 VALUES LESS THAN(2000)(

-□SUBPARTITION s0

-□DATA DIRECTORY='/disk0/data'

-□INDEX DIRECTORY='/disk0/idx',

-□SUBPARTITION s1

-□DATA DIRECTORY='/disk1/data'

-□INDEX DIRECTORY='/disk1/idx'

-□),

-□PARTITION p1 VALUES LESS THAN(2010)(

-□SUBPARTITION s2

-□DATA DIRECTORY='/disk2/data'

-□INDEX DIRECTORY='/disk2/idx',

-□SUBPARTITION s3

-□DATA DIRECTORY='/disk3/data'

-□INDEX DIRECTORY='/disk3/idx'

-□),

-□PARTITION p2 VALUES LESS THAN MAXVALUE(

-□SUBPARTITION s4

-　DATA DIRECTORY='/disk4/data'

-　INDEX DIRECTORY='/disk4/idx',

-　SUBPARTITION s5

-　DATA DIRECTORY='/disk5/data'

-　INDEX DIRECTORY='/disk5/idx'

-　)

-　);

Query OK,0 rows affected(0.02 sec)

---

# 但　InnoDB存储引擎下，建立分区表时并不会按照子分区中的DATA DIRECTORY和INDEX DIRECTORY来创建物理上的文件夹，因此数据和索引文件不会分别存储。

---

mysql　CREATE TABLE ts(a INT,b DATE)engine=innodb

-　PARTITION BY RANGE(YEAR(b))

-　SUBPARTITION BY HASH(TO_DAYS(b))(

-　PARTITION p0 VALUES LESS THAN(2000)(

-　SUBPARTITION s0

-　DATA DIRECTORY='/disk0/data'

-　INDEX DIRECTORY='/disk0/idx',

-　SUBPARTITION s1

-　DATA DIRECTORY='/disk1/data'

-　INDEX DIRECTORY='/disk1/idx'

-　),

-　PARTITION p1 VALUES LESS THAN(2010)(

-　SUBPARTITION s2

-　DATA DIRECTORY='/disk2/data'

-　INDEX DIRECTORY='/disk2/idx',

-　SUBPARTITION s3

-　DATA DIRECTORY='/disk3/data'

-　INDEX DIRECTORY='/disk3/idx'

-□),

-□PARTITION p2 VALUES LESS THAN MAXVALUE(

-□SUBPARTITION s4

-□DATA DIRECTORY='/disk4/data'

-□INDEX DIRECTORY='/disk4/idx',

-□SUBPARTITION s5

-□DATA DIRECTORY='/disk5/data'

-□INDEX DIRECTORY='/disk5/idx'

-□)

-□);

Query OK,0 rows affected(0.02 sec)

mysql□system ls-lh/usr/local/mysql/data/test2/ts*

-rw-rw----1 mysql mysql 8.4K Aug 1 16:24/usr/local/mysql/data/test2/ts.frm

-rw-rw----1 mysql mysql 80 Aug 1 16:24/usr/local/mysql/data/test2/ts.par

-rw-rw----1 mysql mysql 96K Aug 1 16:25/usr/local/mysql/data/test2/ts#P#p0#SP#s0.ibd

-rw-rw----1 mysql mysql 96K Aug 1 16:25/usr/local/mysql/data/test2/ts#P#p0#SP#s1.ibd

-rw-rw----1 mysql mysql 96K Aug 1 16:25/usr/local/mysql/data/test2/ts#P#p1#SP#s2.ibd

-rw-rw----1 mysql mysql 96K Aug 1 16:25/usr/local/mysql/data/test2/ts#P#p1#SP#s3.ibd

-rw-rw----1 mysql mysql 96K Aug 1 16:25/usr/local/mysql/data/test2/ts#P#p2#SP#s4.ibd

-rw-rw----1 mysql mysql 96K Aug 1 16:25/usr/local/mysql/data/test2/ts#P#p2#SP#s5.ibd

# 4.8.4 □□□□□NULL□

MySQL□□□□□□□NULL□□□□□□□□□□□□□□□□□□□□□□□□□□□MYSQL□□□□□
□□□□□□NULL□□□□□□□□□□□□NULL□□□□□MySQL□□□□□□□NULL□□□
ORDER BY□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□NULL□□□□□□□□□
□□□□

□□RANGE□□□□□□□□□□□□□□□NULL□□□□□MySQL□□□□□□□□□□□□□□□□□□□
□□

---

```
mysql□CREATE TABLE t_range(

-□a INT,

-□b INT)ENGINE=InnoDB

-□PARTITION BY RANGE(b)(

-□PARTITION p0 VALUES LESS THAN(10),

-□PARTITION p1 VALUES LESS THAN(20),

-□PARTITION p2 VALUES LESS THAN MAXVALUE

-□);

Query OK,0 rows affected(0.01 sec)
```

---

## □□□□□□□□□□1,1□□□□1,NULL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

```
mysql□INSERT INTO t_range SELECT 1,1;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO t_range SELECT 1,NULL;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql□SELECT*FROM t_range\G;

***************************1.row***************************
```

a:1

b:1

***************************2.row***************************

a:1

b:NULL

2 rows in set(0.00 sec)

mysql⬜SELECT table_name,partition_name,table_rows

-⬜FROM information_schema.PARTITIONS

-⬜WHERE table_schema=DATABASE()AND table_name='t_range'\G;

***************************1.row***************************

table_name:t_range

partition_name:p0

table_rows:2

***************************2.row***************************

table_name:t_range

partition_name:p1

table_rows:0

***************************3.row***************************

table_name:t_range

partition_name:p2

table_rows:0

3 rows in set(0.00 sec)

---

可以看到，两条记录都放入分区p0中了，这是因为对于RANGE分区的插入，NULL值被当作最小值来处理。因此，对于分区p0，存放的是小于等于10的记录，而把NULL值看作小于等于最小的值。

---

mysql⬜ALTER TABLE t_range DROP PARTITION p0;

Query OK,0 rows affected(0.01 sec)

Records:0 Duplicates:0 Warnings:0

mysql：SELECT*FROM t_range;

Empty set(0.00 sec)

---

# 在LIST分区中如果插入NULL值，需要在定义分区时特意指定一个NULL值的分区，否则会报错

---

mysql：CREATE TABLE t_list(

-：a INT,

-：b INT)ENGINE=INNODB

-：PARTITION BY LIST(b)(

-：PARTITION p0 VALUES IN(1,3,5,7,9),

-：PARTITION p1 VALUES IN(0,2,4,6,8)

);

Query OK,0 rows affected(0.00 sec)

mysql：INSERT INTO t_list SELECT 1,NULL;

ERROR 1526(HY000):Table has no partition for value NULL

---

# 在p0分区加入NULL值，就可以正常插入了

---

mysql：CREATE TABLE t_list(

-：a INT,

-：b INT)ENGINE=INNODB

-：PARTITION BY LIST(b)(

-：PARTITION p0 VALUES IN(1,3,5,7,9,NULL),

-：PARTITION p1 VALUES IN(0,2,4,6,8)

);

Query OK,0 rows affected(0.00 sec)

mysql：INSERT INTO t_list SELECT 1,NULL;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql：SELECT table_name,partition_name,table_rows

-　FROM information_schema.PARTITIONS

-　WHERE table_schema=DATABASE()AND table_name='t_list'\G;

***************************1.row***************************

table_name:t_list

partition_name:p0

table_rows:1

***************************2.row***************************

table_name:t_list

partition_name:p1

table_rows:0

2 rows in set(0.00 sec)

---

# HASH、KEY分区对于NULL值的处理和RANGE分区、LIST分区不一样，任何分区函数都会将含有NULL值的记录返回到0分区。

---

mysql　CREATE TABLE t_hash(

-　a INT,

-　b INT)ENGINE=InnoDB

-　PARTITION BY HASH(b)

-　PARTITIONS 4;

Query OK,0 rows affected(0.00 sec)

mysql　INSERT INTO t_hash SELECT 1,0;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql　INSERT INTO t_hash SELECT 1,NULL;

Query OK,1 row affected(0.01 sec)

Records:1 Duplicates:0 Warnings:0

mysql　SELECT table_name,partition_name,table_rows

-　FROM information_schema.PARTITIONS

-　WHERE table_schema=DATABASE()AND table_name='t_hash'\G;

```
***************************1.row***************************

table_name:t_hash

partition_name:p0

table_rows:2

***************************2.row***************************

table_name:t_hash

partition_name:p1

table_rows:0

***************************3.row***************************

table_name:t_hash

partition_name:p2

table_rows:0

***************************4.row***************************

table_name:t_hash

partition_name:p3

table_rows:0

4 rows in set(0.00 sec)
```

# 4.8.5 □□□□□□

□□□□□□□□□□□"□□□□□□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□OLTP□□□□□□□□□□□□Blog□□□□□□□□□□□□□□□□□□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLTP□□□□□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□□□OLTP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLAP□□□□□

□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□1□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Partition Pruning□□□□□

□□□□□OLTP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□10%□□□□□□□□□□□□□□□□□□□□□□□□□□□□□B+□□□□□□□□□□□□□□□□□□□□□□B+□□□2□3□□□□□IO□□□□B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□1000W□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□10□HASH□□□□□□□□□□□□□□□□100W□□□□□□□□□□□□□□□□□□□□SELECT*FROM TABLE WHERE PK=@pk□□□□□□□□□□□□□□□□□□□□□□100W□1000W□□□□□□□□□□□□B+□□□□□□□□□□□□□□□□2□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□1000W□B+□□□□□3□100W□B+□□□□□□2□□□□□□□□□□□□□□□□□□□□□□1□IO□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□SELECT*FROM TABLE WHERE KEY=@key□□□□□□□KEY□□□□□□□□□□□□10□□□□□□□□□□□□□□□□□□□□2□IO□□□□□□□□20□IO□□□□□□□□□□□□□□□□KEY□□□□□□□2□3□IO□

□□□□□□□□□□□Profile□□□□□□□ID□□□□HASH□□□□HASH□□□□□□□□10□□Profile□□□□□1000W□□□□□

---

mysql□CREATE TABLE'Profile'(

-□'id'int(11)NOT NULL AUTO_INCREMENT,

-□'nickname'varchar(20)NOT NULL DEFAULT'',

-□'password'varchar(32)NOT NULL DEFAULT'',

-□'sex'char(1)NOT NULL DEFAULT'',

-□'rdate'date NOT NULL DEFAULT'0000-00-00',

-□PRIMARY KEY('id'),

-□KEY'nickname'('nickname')

-□)ENGINE=InnoDB

-□PARTITION BY HASH(id)

-□PARTITIONS 10;

Query OK,0 rows affected(1.29 sec)

mysql□SELECT COUNT(nickname)FROM Profile;

***************************1.row**************************

count(1):9999248

1 row in set(1 min 24.62 sec)

# 可以看到，HASH分区对数据进行分区的时候，每个分区的数据是大致一样多的

mysql□SELECT table_name,partition_name,table_rows

-□FROM information_schema.PARTITIONS

-□WHERE table_schema=DATABASE()AND table_name='Profile'\G;

***************************1.row**************************

table_name:Profile

partition_name:p0

table_rows:990703

***************************2.row**************************

table_name:Profile

partition_name:p1

table_rows:1086519

***************************3.row**************************

table_name:Profile

partition_name:p2

table_rows:976474

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*4.row\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

table_name:Profile

partition_name:p3

table_rows:986937

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*5.row\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

table_name:Profile

partition_name:p4

table_rows:993667

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*6.row\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

table_name:Profile

partition_name:p5

table_rows:978046

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*7.row\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

table_name:Profile

partition_name:p6

table_rows:990703

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*8.row\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

table_name:Profile

partition_name:p7

table_rows:978639

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*9.row\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

table_name:Profile

partition_name:p8

table_rows:1085334

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*10.row\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

table_name:Profile

partition_name:p9

table_rows:982788

10 rows in set(0.80 sec)

说明　此处的分区标识符是基于HASH函数所计算得到的，它是一个主键，基于ID字段。现在，让我们来看一个相反的例子，若要查找的是一个没有建立索引的字段，会发生什么情况呢？

mysql　EXPLAIN PARTITIONS SELECT*FROM Profile WHERE id=1\G;

***************************1.row***************************

id:1

select_type:SIMPLE

table:Profile

partitions:p1

type:const

possible_keys:PRIMARY

key:PRIMARY

key_len:4

ref:const

rows:1

Extra:

1 row in set(0.00 sec)

## 分区标识符仍然是p1，现在我们来看表Profile中nickname字段的查询结果EXPLAIN PARTITIONS语句的执行结果。

mysql　EXPLAIN PARTITIONS

-　SELECT*FROM Profile WHERE nickname='david'\G;

***************************1.row***************************

id:1

select_type:SIMPLE

table:Profile

partitions:p0,p1,p2,p3,p4,p5,p6,p7,p8,p9

type:ref

possible_keys:nickname

key:nickname

key_len:62

ref:const

rows:10

Extra:Using where

1 row in set(0.00 sec)

---

## 此时看看MySQL数据库是否增加了，这对于原本看上去比较神秘的数据库也就不再神秘了。

---

mysql□SELECT*FROM Profile WHERE nickname='david'\G;

***************************1.row**************************

id:5566

nickname:david

password:3e35d1025659d07ae28e0069ec51ab92

sex:M

rdate:2003-09-20

1 row in set(1.05 sec)

---

## 可以看到，上述查询用了1.05秒，这是因为数据库刚刚启动，缓冲池中还没有缓存任何数据，需要IO读取大约20～30次的磁盘操作。当再次执行上述的查询语句时，执行的SQL语句时间为0.26秒。

## 由此可以发现，InnoDB存储引擎是基于OLTP应用设计的，而大多数情况下的查询都是基于主键进行的，对于OLTP应用来说，用户一般只需要通过索引查询少量的数据，而不是大量的数据。

# 4.8.6 交换分区与子分区数据

MySQL 5.6支持使用ALTER TABLE…EXCHANGE PARTITION语句来进行分区或子分区与表之间的数据交换。如果分区或子分区中的数据比较多，可以把这些数据交换到一张表中去，能够更快地删除分区中的数据。反之也可以把表中的数据写入到一个分区中去。

使用ALTER TABLE…EXCHANGE PARTITION语句时需要注意以下内容：

❑用来交换的表的结构必须和分区表的结构一致，即表的结构相同。

❑用来交换的表不能包含任何分区或子分区。

❑用来交换的表中不能含有外键，或者其他表引用的外键。

❑用户除了需要ALTER、INSERT、CREATE权限外，还需要DROP的权限。

❑在交换过程中不会触发交换表和被交换表上的触发器。

❑AUTO_INCREMENT列将被重置。

下面举例说明，首先创建一个含有RANGE分区的表e，示例代码如下所示。

---

```
CREATE TABLE e(

id INT NOT NULL,

fname VARCHAR(30),

lname VARCHAR(30)

)

PARTITION BY RANGE(id)(

PARTITION p0 VALUES LESS THAN(50),

PARTITION p1 VALUES LESS THAN(100),

PARTITION p2 VALUES LESS THAN(150),

PARTITION p3 VALUES LESS THAN(MAXVALUE)
```

```
);

INSERT INTO e VALUES

(1669,"Jim","Smith"),

(337,"Mary","Jones"),

(16,"Frank","White"),

(2005,"Linda","Black");
```

# 如果你想创建一个e2表，让e2是一个非分区表，其数据和表e相同，但是你不想e2也是一个分区表

```
mysql>CREATE TABLE e2 LIKE e;

Query OK,0 rows affected(1.34 sec)

mysql>ALTER TABLE e2 REMOVE PARTITIONING;

Query OK,0 rows affected(0.90 sec)

Records:0 Duplicates:0 Warnings:0
```

# 获取分区表的信息，分区表的记录数等

```
mysql>SELECT PARTITION_NAME,TABLE_ROWS

->FROM INFORMATION_SCHEMA.PARTITIONS

->WHERE TABLE_NAME='e';

+----------------+------------+

|PARTITION_NAME|TABLE_ROWS|

+----------------+------------+

|p0|1|

|p1|0|

|p2|0|

|p3|3|

+----------------+------------+

4 rows in set(0.00 sec)
```

## 现在把e2表中的数据加载到原始分区表e的原来p0分区中。现在看看e2表。

mysql＞ALTER TABLE e EXCHANGE PARTITION p0 WITH TABLE e2;

Query OK,0 rows affected(0.28 sec)

## 现在可以看到表e中满足分区规则的p0分区中已经有数据了。

mysql＞SELECT PARTITION_NAME,TABLE_ROWS

-＞FROM INFORMATION_SCHEMA.PARTITIONS

-＞WHERE TABLE_NAME='e';

+----------------+------------+

|PARTITION_NAME|TABLE_ROWS|

+----------------+------------+

|p0|0|

|p1|0|

|p2|0|

|p3|3|

+----------------+------------+

4 rows in set(0.00 sec)

## 下面可以看到原来e2表中的数据已经没有了。

mysql＞SELECT*FROM e2;

+----+-------+-------+

|id|fname|lname|

+----+-------+-------+

|16|Frank|White|

+----+-------+-------+

1 row in set(0.00 sec)

# 4.9　小结

本章主要介绍了表的结构、InnoDB存储引擎表的物理存储方式及数据的组织方式。通过InnoDB存储引擎表的物理存储方式，读者可以更加深入地理解表、页、行的概念。同时，本章还介绍了分区的概念，MySQL数据库支持分区功能，通过分区可以将一张表的数据按照InnoDB存储引擎表的物理存储方式进行划分。需要注意的是，MySQL数据库的分区是局部分区索引，一个分区中既存放了数据又存放了索引，因此在使用时需要注意其性能问题。

本章还介绍了当前MySQL支持的分区RANGE、LIST、HASH、KEY、COLUMNS类型，一般来说，使用HASH、KEY这类分区方式需要注意其使用的场景。对于OLTP的应用，建议不要对其进行分区，否则可能会带来性能问题。

# 第5章　优化服务器设置

在本章中，我们将讨论如何为运行得更稳定、更持久创建一个合理的配置文件。尽管有大量的选项，但从根本上来说，只有很少的一部分是真正重要的。

如果因为不恰当的配置而陷入麻烦——我们看到过很多因为试图对一些参数进行优化而适得其反的例子。在前面的案例中，我们曾经见过由于有人根据一些似是而非的教条调整了SQL服务器，导致生产服务器崩溃的例子，不幸的是这样的DBA有很多。大部分DBA不敢随意调整服务器的参数，因为他们并不理解这些参数到底有什么作用。结果导致有些SQL服务器的配置实际上和默认值差别很大。有时候这么做是没有问题的，但有时候确实非常危险。举个例子，我们曾经见过一个MySQL，通过iostat发现磁盘的利用率一直是100%，但过了一会儿就恢复正常了，原来是有人将服务器上的缓冲刷到磁盘的比例改成了20%，这在大部分时候都工作得很好，只不过偶尔会带来问题。

另外一个常见的错误是InnoDB日志文件的大小，我们看到过很多因为日志文件太小而导致服务器性能很差的案例。但在某些情况下，MySQL的默认配置也是非常糟糕的，有些配置需要根据硬件来修改。

# 5.1 InnoDB □□□□□□□□□

InnoDB□□□□□□□□□□□□□□□□□□□□□□

□B+□□□□

□□□□□□

□□□□□□

□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□B+□□□□□□□□□□□□□□□□□□□□□Key Value□□□□□□□□□□

□□ B+□□□B□□□□□□□□binary□□□□□□□□□□balance□□□□□B+□□□□□□□□□□□□□□□□□□□□□□B+□□□□□□□□□□

□□□□□□□DBA□□□□□□□□□B+□□□□□□□□□□□□□□□□□□□□□□□□□B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 5.2　随机读取与顺序

B+树索引的特点是能够高效地找到一条记录或者一个范围的记录，而哈希表结构只能用于查找单条记录。本节将主要介绍B+树索引的相关知识。

## 5.2.1　二分查找法

二分查找法（binary search）也称折半查找法，用来查找一组有序的记录数组中的某一记录，其基本思想是：将记录按有序化（递增或递减）排列，在查找过程中采用跳跃式方式查找，即先以有序数列的中点位置为比较对象，如果要找的元素值小于该中点元素，则将待查序列缩小为左半部分，否则为右半部分。通过一次比较，将查找区间缩小一半。

如有5、10、19、21、31、37、42、48、50、55这10个数，现要从这10个数中查找48这条记录，其查找过程如图5-1所示。



图　5-1　二分查找法

如图5-1所示，将第3列指向第48号页面的页表项的访问位设置为8，将该访问位较小的列向前移动。第5号页面的访问位最小为1，所以将其放在第4列。将该列的访问位设置为10。访问位的平均值为（1+2+3+4+5+6+7+8+9+10）/10＝5.5，访问位的方差为（4+3+2+4+3+1+4+3+2+3）/10＝2.9。由于第三列的访问位为10，所以将其放在第4列。

由于以上方法需要较多的存储空间和较长的处理时间，因此在1946年提出的方法中，只使用一个访问位。在1962年提出的方法中，将访问位分为两个部分，即Page Directory部分和页表部分。访问位放在页表中，与访问位对应的Page Directory中的访问位设置为零。

# 5.2.2 索引文件和数据行的存储

首先，B+树索引只是实现了索引的数据结构，而B+树索引的每个叶子节点又是如何与B树索引和数据行建立起对应关系的呢？我们通过一个简单的例子来讲解两者之间的关系，如图5-2所示的二叉搜索树。



图 5-2 二叉搜索树

图5-2所示二叉搜索树中的每个节点都是数据行，每个数据行都有一个指向其左孩子和右孩子的指针。图5-2所示的二叉搜索树共有六个数据行：2、3、5、6、7、8。

在图5-2所示的二叉搜索树中，如果想查找数据行5，首先会与根节点数据行6（6是比5大的最小的6）比较，3比5大，5比3大，于是向右转，与数据行3比较，数据行有2、3、5、6、7、8三个节点，于是与数据行3比较，当只有一个数据行8的时候，数据行3的时候，直接与6进行比较，接下来的比较次数非常接近，节点查询的平均比较次数是（1+2+3+4+5+6）/6＝3.3，节点查询的平均比较次数是（3+3+3+2+2+1）/6＝2.3，查询效率明显上升了，比较次数明显减少。

如果想查找从小到大排序的数据行2、3、5、6、7、8，那么如何进行查找呢，如图5-3所示，即为排序后的链表。

（5-3）。该二叉树的平均查找长度为（1+2+3+4+5+5）/6＝3.16。从查找效率看，显然该二叉排序树是不"平衡"的，我们更希望得到如左图那样比较"矮"的树，因为它的查找效率更高。这就引出了我们下面将要学习的——平衡二叉树的概念（AVL树）。



图 5-3　不平衡的二叉排序树示意图

之所以会产生这种不平衡的情况，是因为我们在插入节点的时候，没有考虑到每插入一个节点后（如图5-3）对整个树结构的影响。如图5-2所示，如果我们在插入的过程中，采取一定的"平衡"措施，尽量使得任何时候树的结构都是平衡的，那么就能够有效提高查找效率，从而避免上面那种不平衡的情况发生。

平衡二叉排序树在进行插入操作的时候，每插入一个节点，都要保证任何节点到其两个子树的高度差不超过1。利用这个性质可以有效地控制树的高度（如图5-2所示），从而避免产生普通的二叉排序树（如图5-3所示那样）的不平衡情况。

附录5-4：调查问卷表

插入新值9



左旋以保证平衡

图 5-4　金属层有9个布线通道的单元

为了提高电路的集成度，金属层的最大利用率就显得十分重要。如图5-5所示，

平衡二叉树

插入新键值3

右旋一次

再左旋一次

## □ 5-5 □□□□□□□□□□□□

□5-4□□5-5□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 5.3  B+树

B+树索引是目前关系型数据库系统中查找最为常用和最为有效的索引。B+树的B不是代表二叉（binary），而是代表平衡（balance），因为B+树是从最早的平衡二叉树ISAM树演化而来，但是B+树不是一个二叉树。MyISAM引擎使用的就是这个索引结构。下面我们将从最早的平衡二叉树开始，一步一步地演绎B+树的由来。

B+树是为磁盘或其他直接存取辅助设备设计的一种平衡查找树。在B+树中，所有记录节点都是按键值的大小顺序存放在同一层的叶子节点上，由各叶子节点指针进行连接。先来看一个B+树，其高度为2，每页可存放4条记录，扇出（fan out）为5，如图5-6所示。



图  5-6    高度为2的B+树

从图5-6可以看出，所有记录都在叶子节点上，并且是顺序存放的，如果用户从最左边的叶子节点开始顺序遍历，可以得到所有键值的顺序排序：5、10、15、20、25、30、50、55、60、65、75、80、

## 5.3.1　B+树基本概念

B+树是应文件系统所需而出现的一种数据结构，它是数据库管理系统的基础。B+树是平衡树，其所有叶子节点在同一层上的基本结构如图5-1所示。

表 5-1  B+ 树插入的 3 种情况

| Leaf Page 满 | Index Page 满 | 操作 |
|---|---|---|
| No | No | 直接将记录插入到叶子节点 |
| Yes | No | 1）拆分 Leaf Page<br>2）将中间的节点放入到 Index Page 中<br>3）小于中间节点的记录放左边<br>4）大于或等于中间节点的记录放右边 |
| Yes | Yes | 1）拆分 Leaf Page<br>2）小于中间节点的记录放左边<br>3）大于或等于中间节点的记录放右边<br>4）拆分 Index Page<br>5）小于中间节点的记录放左边<br>6）大于中间节点的记录放右边<br>7）中间节点放入上一层 Index Page |

为了保持插入后整棵B+树的平衡，在对上述3种情况进行插入操作后，需要对B+树进行处理。下面我们用一个例子来分析B+树的插入。例如，对于图5-6所示的B+树，若用户插入28这个键值，则由于其满足Leaf Page、Index Page都不满的情况，所以直接进行插入即可，完成后如图5-7。

图 5-7 插入键值28

接着插入70，这时原先的Leaf Page已满，需要进行Index Page的拆分操作，这时图5-1中所示的页中的Leaf Page的键值分别为55、55、60、65、70，拆分点为60，很显然拆分后键值见图5-8。

图 5-8 插入键值70

接着插入键值95，这时同样需要设置为页面进行分裂，和5-6页5-7插入所做的操作类似，

而插入键值95后又需要对5-1结构进行分裂，此时对Leaf Page和Index Page都进行了

分裂，分裂后的结果如图5-9所示。

图 5-9 插入关键字95

一次插入操作有可能导致B+树的一个叶结点分裂，进而导致其父结点分裂，父结点的分裂又可能导致祖父结点的分裂（split），所以说在B+树中插入一个关键字，有可能引起整个树的结构的改变，但结构的改变只会从叶结点开始逐层向上传播。在B+树中删除一个关键字也有同样的情况发生，这里涉及旋转（Rotation）等操作。

叶子节点（Leaf Page）存放数据，所以相同数据量的情况下，没有存放索引键的B+树，比有存放索引键的二叉树更加矮胖，所需要从磁盘读取的页面数据会更少。另外，图5-7中二叉树的高度为70，而B+树的高度只有两层，只需要磁盘读取三次即可，如图5-10所示的索引。



图 5-10    B+树索引示意图

在图5-10中，可以看出，相同数据量下的B+树索引要比二叉树矮胖得多，该B+树索引的高度为2。

## 5.3.2　B+树的删除操作

B+树使用填充因子（fill factor）来控制树的删除变化，50%是填充因子可设的最小值。B+树的删除操作同样必须保证删除后叶子节点中的记录依然排序，同插入一样，B+树的删除操作同样需要考虑三种情况，与插入不同的是，删除根据填充因子的变化来衡量。

表 5-2　B+ 树删除操作的三种情况

| 叶子节点小于填充因子 | 中间节点小于填充因子 | 操作 |
| --- | --- | --- |
| No | No | 直接将记录从叶子节点删除，如果该节点还是 Index Page 的节点，用该节点的右节点代替 |
| Yes | No | 合并叶子节点和它的兄弟节点，同时更新 Index Page |
| Yes | Yes | 1）合并叶子节点和它的兄弟节点<br>2）更新 Index Page<br>3）合并 Index Page 和它的兄弟节点 |

接着用图5-9的B+树来进行删除操作。首先删除键值为70的这条记录，该记录符合表5-2讨论的第一种情况，删除后可得到图5-11。

图 5-11  插入数据70

25　　　　　　　5-2　　　　　　　　　　　IndexPage　　　　　　　　Leaf Page　　25　　　　25　　　　　　28　　　Page Index　　　　　　5-12

图 5-12　删除结点25

如果我们再插入60这个键值，Leaf Page就必须为60分裂，由于Fill Factor是50%，所以分裂后的数据重新分配，并且在上层Index Page分裂出一个新的Index Page，整个数据分布如图5-13。



图 5-13 键值插入60

# 5.4 B+树索引

B+树索引的本质就是B+树在数据库中的实现。但是B+索引在数据库中有一个特点是高扇出性，因此在数据库中，B+树的高度一般都在2～4层，这也就是说查找某一键值的行记录时最多只需要2～4次IO，这倒不错。因为当前一般的机械磁盘每秒至少可以做100次IO，2～4次的IO意味着查询时间只需0.02～0.04秒。

数据库中的B+树索引可以分为聚集索引（clustered inex）和辅助索引（secondary index）[1]，但是不管是聚集还是辅助的索引，其内部都是B+树的，即高度平衡的，叶子节点存放着所有的数据。聚集索引与辅助索引不同的是，叶子节点存放的是否是一整行的信息。

## 5.4.1 聚集索引

之前已经介绍过了，InnoDB存储引擎表是索引组织表，即表中数据按照主键顺序存放。而聚集索引（clustered index）就是按照每张表的主键构造一棵B+树，同时叶子节点中存放的即为整张表的行记录数据，也将聚集索引的叶子节点称为数据页。聚集索引的这个特性决定了索引组织表中数据也是索引的一部分。同B+树数据结构一样，每个数据页都通过一个双向链表来进行链接。

由于实际的数据页只能按照一棵B+树进行排序，因此每张表只能拥有一个聚集索引。在多数情况下，查询优化器倾向于采用聚集索引。因为聚集索引能够在B+树索引的叶子节点上直接找到数据。此外，由于定义了数据的逻辑顺序，聚集索引能够特别快地访问针对范围值的查询。查询优化器能够快速发现某一段范围的数据页需要扫描。

下面通过例子来分析聚集索引中数据是如何存储的。首先创建一张表t：

---

```
CREATE TABLE t（

a INT NOT NULL,

b VARCHAR(8000),

c INT NOT NULL,

PRIMARY KEY(a),

KEY idx_c(c)

)ENGINE=INNODB;

INSERT INTO t SELECT 1,REPEAT('a',7000),-1;

INSERT INTO t SELECT 2,REPEAT('a',7000),-2;
```

INSERT INTO t SELECT 3,REPEAT('a',7000),-3;

INSERT INTO t SELECT 4,REPEAT('a',7000),-4;

---

## 这次再来观察表t，由于b列最大7000字节，这次插入4条记录后，数据页已经产生了分裂，通过py_innodb_page_info工具可以发现这次分裂：

---

[root@nineyou0-43 data]#py_innodb_page_info.py-v mytest/t.ibd

page offset 00000000,page type<File Space Header>

page offset 00000001,page type<Insert Buffer Bitmap>

page offset 00000002,page type<File Segment inode>

page offset 00000003,page type<B-tree Node>,page level<0001>

page offset 00000004,page type<B-tree Node>,page level<0000>

page offset 00000005,page type<B-tree Node>,page level<0000>

page offset 00000006,page type<B-tree Node>,page level<0000>

page offset 00000000,page type<Freshly Allocated Page>

Total number of page:8:

Freshly Allocated Page:1

Insert Buffer Bitmap:1

File Space Header:1

B-tree Node:4

File Segment inode:1

---

## page level为0000的即为叶子节点，而数据页就是存放在叶子节点上的，那么根据上面的分析，我们可以通过page level为0001的非叶子节点来观察B+树内部2层结构，B+树内部结构通过hexdump来观察，可以在其中看到如下的数据：

---

0000c000 c2 33 62 95 00 00 00 03 ff ff ff ff ff ff ff ff|.3b............|

0000c010 00 00 00 0a b6 8c ce 57 45 bf 00 00 00 00 00 00|.......WE.......|

0000c020 00 00 00 00 00 f9 00 02 00 a2 80 05 00 00 00 00|...............|

0000c030 00 9a 00 02 00 02 00 03 00 00 00 00 00 00 00 00|...............|

```
0000c040 00 01 00 00 00 00 00 00 01 e2 00 00 00 f9 00 00|...............|

0000c050 00 02 00 f2 00 00 00 f9 00 00 00 02 00 32 01 00|.............2..|

0000c060 02 00 1b 69 6e 66 69 6d 75 6d 00 04 00 0b 00 00|...infimum......|

0000c070 73 75 70 72 65 6d 75 6d 00 10 00 11 00 0e 80 00|supremum........|

0000c080 00 01 00 00 00 04 00 00 00 19 00 0e 80 00 00 02|...............|

0000c090 00 00 00 05 00 00 00 21 ff d6 80 00 00 04 00 00|.......!........|

0000c0a0 00 06 00 00 00 00 00 00 00 00 00 00 00 00 00 00|...............|

0000c0b0 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00|...............|

0000c0c0 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00|...............|

......

0000fff0 00 00 00 00 00 70 00 63 73 d8 52 3a b6 8c ce 57|.....p.cs.R:...W|
```

接下来我们不再分析Page Directory部分。我们从00 63这个地址开始向前看，这是第一个记录的Recorder Header，紧接着0xc063的数据是69 6e 66 69 6d 75 6d 00，也就是infimum。再向前分析5字节01 00 02 00 1b这是Recorder Header，一共有4个字8个字节1个字节。这里我们不再具体分析，只需要知道InnoDB的Page Directory的值指向的是infimum，而这个记录的Recorder Header的最后两个字节00 1b，表示下一个记录的地址c063+1b=c07e，这里的数据是80 00 00 01，也就是第一行1这条记录，因为是INT无符号类型的数据，所以0x80 00 00 01表示的是0x0001。而80 00 00 01的下个00 00 00 04这又是下一个记录地址，我们可以看到80 00 00 02、80 00 00 04这也就是我们的每行主键的值。

我们也可以用类似的方式，将这个数据完全解析。用这个方法完全读懂这个数据以后我们就可以做一些实验了，接下来我们通过实验来证实下前面讨论过的5-14的内容

图 5-14 B+树结构

如果表上有其他索引，那么同样可以使用相似的方法进行测试，这里就不再赘述。图5-14给出的是一个简单的示意图，读者可以在此基础上进行更深入的分析和研究，以便更好地理解。

当数据库处理查询时，它会读取所有需要的行，然后根据指定的排序规则对结果进行排序。如果可以使用索引来进行排序，那么排序过程将非常快速，否则就会很慢。

下面是一个没有使用索引排序的示例，虽然它使用了主键索引，但实际上是按照主键的顺序来读取数据的。数据库会先读取前10条记录（利用B+树索引快速定位），然后直接返回这些记录。下面展示了前10条记录的具体内容，EXPLAIN的输出结果如下：

---

mysql›EXPLAIN

-›SELECT*FROM Profile ORDER BY id LIMIT 10\G;

***************************1.row***************************

id:1

select_type:SIMPLE

table:Profile

type:index

possible_keys:NULL

key:PRIMARY

key_len:4

ref:NULL

rows:10

Extra:

1 row in set(0.00 sec)

---

在这个查询中，由于ORDER BY子句和查询使用了相同的索引，所以就不需要额外的filesort操作，这样查询速度会非常快。

下面是一个典型的range query（范围查询）示例，它同样使用了主键索引。从查询优化器的角度来看，它能够快速定位到满足条件的记录并返回结果。

---

mysql›EXPLAIN

-›SELECT*FROM Profile

-›WHERE id›10 AND id‹10000\G;

***************************1.row***************************

id:1

select_type:SIMPLE

table:Profile

type:range

possible_keys:PRIMARY

key:PRIMARY

key_len:4

ref:NULL

rows:14868

Extra:Using where

1 row in set（0.01 sec）

---

□□EXPLAIN□□□□MySQL□□□□□□□□□□execute plan□□□□□rows□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□rows□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□9946□□□□□

---

mysql□SELECT COUNT(*)from Profile

-□WHERE id□10 AND id□10000;

***************************1.row*************************

COUNT(1):9946

1 row in set(0.00 sec)

---

[1]_□□□□□□□□□□□□□□□□□□□non-clustered index□□□

## 5.4.2 　辅助索引

辅助索引（也称Secondary Index）也叫非聚集索引。辅助索引在叶子节点中并不会存放所有数据，辅助索引叶子节点存储的是键值，以及书签（bookmark）。这个书签就是对应的InnoDB的聚集索引的键。这样的策略带来了InnoDB获取真实数据比聚集索引慢，InnoDB需要获取两次索引。下面就让我们看看如图5-15所示的InnoDB辅助索引的实现方式，进行进一步学习。

Secondary Index

Secondary Index

Clustered Index

# 图 5-15 独立表空间的存储结构

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□3□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□3□□□□□□□□□□□□□□□□□□□□3□□□□□□□□□□□□□□□□3□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□6□□□IO□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□Microsoft SQL Server□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□MyISAM□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□NOT NULL□□□□□□□□□□□□□□□□□Row Identifiedr□RID□□□□□□□"□□□□□□□□□□□"□□□□□□□□□□□□□□□□

□□Microsoft SQL Server□□□□DBA□□□□□□□□□□□□□□□Microsoft SQL Server□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLAP□On-Line Analytical Processing□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□read ahead□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□MyISAM□□□□□□□□□□□□□□□□□□□It all depends□

□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□t□□□□□□□□□t□□□□□□□c□□□□□c□□□□□□□□□□

---

```
mysql□ALTER TABLE t

-□ADD c INT NOT NULL;

Query OK,4 rows affected(0.24 sec)

Records:4 Duplicates:0 Warnings:0

mysql□UPDATE t SET c=0-a;

Query OK,4 rows affected(0.04 sec)
```

Rows matched:4 Changed:4 Warnings:0

mysql>ALTER TALBE t ADDKEY idx_c(c);

Query OK,4 rows affected(0.28 sec)

Records:4 Duplicates:0 Warnings:0

mysql>SHOW INDEX FROM t\G;

***************************1.row***************************

Table:t

Non_unique:0

Key_name:PRIMARY

Seq_in_index:1

Column_name:a

Collation:A

Cardinality:2

Sub_part:NULL

Packed:NULL

Null:

Index_type:BTREE

Comment:

***************************2.row***************************

Table:t

Non_unique:1

Key_name:idx_c

Seq_in_index:1

Column_name:c

Collation:A

Cardinality:2

Sub_part:NULL

Packed:NULL

Null:

Index_type:BTREE

Comment:

2 rows in set（0.00 sec）

mysql）select a,c from t;

+---+----+

|a|c|

+---+----+

|4|-4|

|3|-3|

|2|-2|

|1|-1|

+---+----+

4 rows in set（0.00 sec）

---

# 利用工具py_innodb_page_info分析表空间结构

---

[root@nineyou0-43 mytest]#py_innodb_page_info.py-v t.ibd

page offset 00000000,page type（File Space Header）

page offset 00000001,page type（Insert Buffer Bitmap）

page offset 00000002,page type（File Segment inode）

page offset 00000003,page type（B-tree Node），page level（0001）

page offset 00000004,page type（B-tree Node），page level（0000）

page offset 00000005,page type（B-tree Node），page level（0000）

page offset 00000006,page type（B-tree Node），page level（0000）

page offset 00000007,page type（B-tree Node），page level（0000）

page offset 00000000,page type（Freshly Allocated Page）

Total number of page:9:

Freshly Allocated Page:1

Insert Buffer Bitmap:1

File Space Header:1

B-tree Node:5

当前页面实际上就是该索引页的第一页（page offset为4的一个页面）的内容，其对应的hexdump如下所示：

```
00010000  b9 aa 8e d0 00 00 00 04 ff ff ff ff ff ff ff ff|................|
00010010  00 00 00 0a ec ea 4e 27 45 bf 00 00 00 00 00 00|......N'E.......|
00010020  00 00 00 00 01 02 00 02 00 ac 80 06 00 00 00 00|................|
00010030  00 a4 00 01 00 03 00 04 00 00 00 00 00 52 d4 8b|.............R..|
00010040  00 00 00 00 00 00 00 00 01 f2 00 00 01 02 00 00|................|
00010050  00 02 02 72 00 00 01 02 00 00 00 02 01 b2 01 00|...r............|
00010060  02 00 41 69 6e 66 69 6d 75 6d 00 05 00 0b 00 00|..Ainfimum......|
00010070  73 75 70 72 65 6d 75 6d 00 00 10 ff f3 7f ff ff|supremum........|
00010080  ff 80 00 00 01 00 00 18 ff f3 7f ff ff fe 80 00|................|
00010090  00 02 00 00 20 ff f3 7f ff ff fd 80 00 00 03 00|..............|
000100a0  00 28 ff f3 7f ff ff fc 80 00 00 04 00 00 00 00|.□.............|
……
00013ff0  00 00 00 00 00 70 00 63 f3 46 77 f2 ec ea 4e 27|.....p.c.Fw...N'|
```

可以看到，4号页面类型值（c）为4（即索引页），且其目录项数值和用户记录值正好是我们图5-16中分析的值。

辅助索引idx_c

Page Offset:00 04

| Key:7f ff ff ff | Key:7f ff ff fe | Key:7f ff ff fd | Key:7f ff ff fc |
|---|---|---|---|
| Pointer:80 00 00 01 | Pointer:80 00 00 02 | Pointer:80 00 00 03 | Pointer:80 00 00 04 |

聚集索引|Page Offset:00 03

| Key:80 00 00 01 | Key:80 00 00 02 | Key:80 00 00 04 | |
|---|---|---|---|
| Pointer:00 05 | Pointer:00 06 | Pointer:00 07 | |

Page Offset:00 05

1,repeat（'a',7000），
–1

Page Offset:00 06

2,repeat（'a',7000),–2

3,repeat（'a',7000），–3

Page Offset:00 07

4,repeat（'a',7000),–4

## 表 5-16 中的标志位

由5-16可知，在t时刻，变量idx_c的不同取值对应了不同的标志位含义。当c的值为一个很大的整数时，其十六进制表示为-1（7f ff ff ff），对应的标志位为7（0111）。

当标志位为0时，表示□□□□□□□□；当标志位为1时，表示-1。

# 5.4.3 B+树的旋转操作

在5.3节我们讲述B+树页分裂时有意回避了一种场景，即B+树的旋转。其实这是为了和5.3节讲述的分裂操作进行区分，B+树的旋转发生在插入操作时。

B+树的旋转发生在插入的记录是页的最大值且页空间已满，需要进行页分裂操作时：

1、2、3、4、5、6、7、8、9

此时，由于插入的记录是最大值10，那么其实也就没有必要像之前5.3.1节中介绍的那样，分裂为各含有5条记录的两个页，即split record为6。可以将其分为：

P1：1、2、3、4

P2：5、6、7、8、9、10

可以看到在新插入的页P2中，空出了更多的空间来存放之后插入的记录。然而，在页P2中到底选择哪条记录作为InnoDB存储引擎的Page Header中，有着以下三个部分用来保存插入的顺序信息：

□PAGE_LAST_INSERT

□PAGE_DIRECTION

□PAGE_N_DIRECTION

通过这些信息，InnoDB存储引擎可以决定是向左还是向右进行旋转，同时决定了旋转点记录为哪一个。还是以上述插入为例，若插入的值是5，则根据变量cursor进行判断，InnoDB存储引擎将旋转点记录设为倒数第二条记录。而在插入值为3时，根据之前的信息判断，旋转点记录为倒数第一条。

旋转操作可以很大程度上地减少页的使用，以提高B+树的空间使用率。图5-17所示，

图 5-17　游标记录和拆分记录

在5-17中，我们可以很明显地看出，3者之中split record的位置是最靠右的。有关其具体实现，图5-18给出了代码：

图 5-18 记录过多导致页分裂的情况

如图图5-19所示，页中插入记录后，由于记录过大导致页需要进行分裂，分裂后需要进行调整的情况。

Page

cursor record

split record

record to be insert

After Split

Page

Right Page

cursor record

split record

record to be insert

图 5-19　某工程监理工作流程

# 5.4.4 B+树索引的管理

## 1.索引管理

索引管理该操作同样可以分为两种方法，一种是ALTER TABLE，另一种是CREATE/DROP INDEX。通过ALTER TABLE创建索引的语法为：

---

```
ALTER TABLE tbl_name

|ADD{INDEX|KEY}[index_name]

[index_type](index_col_name,...)[index_option]...

ALTER TABLE tbl_name

DROP PRIMARY KEY

|DROP{INDEX|KEY}index_name
```

---

## CREATE/DROP INDEX的语法同样很简单：

---

```
CREATE[UNIQUE]INDEX index_name

[index_type]

ON tbl_name(index_col_name,...)

DROP INDEX index_name ON tbl_name
```

---

用户可以设置对整个列的数据进行索引，也可以只索引一个列的开头部分数据，如前面创建的表t，列b为varchar（8000），但是在创建索引时只索引了前100个字段，如：

---

```
mysql＞ALTER TABLE t

-＞ADD KEY idx_b(b(100));

Query OK,4 rows affected(0.32 sec)

Records:4 Duplicates:0 Warnings:0
```

---

为了方便读者理解联合索引的内部结构，通过SHOW INDEX命令可以发现，表t中有两个索引，分别是列a、c组成的索引idx_a_c和列

---

mysql□ALTER TABLE t

-□ADD KEY idx_a_c(a,c);

Query OK,4 rows affected(0.31 sec)

Records:4 Duplicates:0 Warnings:0

mysql□SHOW INDEX FROM t\G;

***************************1.row***************************

Table:t

Non_unique:0

Key_name:PRIMARY

Seq_in_index:1

Column_name:a

Collation:A

Cardinality:2

Sub_part:NULL

Packed:NULL

Null:

Index_type:BTREE

Comment:

***************************2.row***************************

Table:t

Non_unique:1

Key_name:idx_b

Seq_in_index:1

Column_name:b

Collation:A

Cardinality:2

Sub_part:100

Packed:NULL

Null:YES

Index_type:BTREE

Comment:

***************************3.row***************************

Table:t

Non_unique:1

Key_name:idx_a_c

Seq_in_index:1

Column_name:a

Collation:A

Cardinality:2

Sub_part:NULL

Packed:NULL

Null:

Index_type:BTREE

Comment:

***************************4.row***************************

Table:t

Non_unique:1

Key_name:idx_a_c

Seq_in_index:2

Column_name:c

Collation:A

Cardinality:2

Sub_part:NULL

Packed:NULL

Null:

Index_type:BTREE

Comment:

```
***************************5.row***************************

Table:t

Non_unique:1

Key_name:idx_c

Seq_in_index:1

Column_name:c

Collation:A

Cardinality:2

Sub_part:NULL

Packed:NULL

Null:

Index_type:BTREE

Comment:

5 rows in set（0.00 sec）
```

_____

通过上面SHOW INDEX FROM命令可以发现，表t中有4个索引，分别为主键，列c上的辅助索引，列b上及列a和c上的组合索引。接着分析通过命令SHOW INDEX展现结果中的各列含义。

☐Table：索引所在的表名。

☐Non_unique：非唯一的索引，可以看到primary key是0，因为必须是唯一的。

☐Key_name：索引的名字，用户可以通过这个名字来执行DROP INDEX。

☐Seq_in_index：索引中该列的位置，如果看组合索引idx_a_c就比较直观了。

☐Column_name：索引列的名称。

☐Collation：列以什么方式存储在索引中。可以是A或NULL。B+树索引总是A，即排序的。如果使用了Heap存储引擎，并且建立了Hash索引，这里就会显示NULL了。因为Hash根据Hash桶存放索引数据，而不是对数据进行排序。

☐Cardinality：非常关键的值，表示索引中唯一值的数目的估计值。Cardinality表的行数应尽可能接近1，如果非常小，那么用户需要考虑是否可以删除此索引。

☐Sub_part：是否是列的部分被索引。如果索引idx_b，这里显示的值为100，表示只对b列的前100字符进行索引。如果索引整个列，则该字段为NULL。

☐Packed：关键字如何被压缩。如果没有被压缩，则为NULL。

☐Null：是否索引的列含有NULL值。可以看到idx_b，这里值为Yes，这是因为定义的列b允许NULL值。

☐Index_type：索引的类型。InnoDB存储引擎只支持B+树索引，所以这里显示的都是BTREE。

☐Comment：注释。

Cardinality值非常关键，优化器会根据这个值来判断是否使用这个索引。但是这个值并不是实时更新的，即并非每次索引的更新都会更新该值，因为这样代价太大了。因此这个值是不太准确的，只是一个大概的值。上面显示的结果中，Cardinality为2，这是因为我们插入了4条记录，但是其中有4条记录为空。如果需要更新索引Cardinality的信息，可以使用ANALYZE TABLE命令，如下：

---

mysql〉analyze table t\G;

***************************1.row***************************

Table:mytest.t

Op:analyze

Msg_type:status

Msg_text:OK

1 row in set(0.01 sec)

mysql〉show index from t\G;

***************************1.row***************************

Table:t

Non_unique:0

Key_name:PRIMARY

Seq_in_index:1

Column_name:a

```
        Collation:A

        Cardinality:4

        Sub_part:NULL

        Packed:NULL

        Null:

        Index_type:BTREE

        Comment:

        ……
```

可以看到Cardinality值非常得不准确，通过命令可以发现此时需要通过ANALYZE TABLE命令来统计索引的信息，该命令执行后使用MySQL数据库对于Cardinality值进行了统计，然后再执行一次，观察到的值就比较精确了。

```
mysql>show index from Profile\G;

***************************1.row***************************

        Table:Profile

        Non_unique:0

        Key_name:UserName

        Seq_in_index:1

        Column_name:username

        Collation:A

        Cardinality:NULL

        Sub_part:NULL

        Packed:NULL

        Null:

        Index_type:BTREE

        Comment:
```

Cardinality为NULL，在某些情况下可能会发生索引建立了却没有用到的情况，或者对于两条基本一样的EXPLAIN语句执行了多次后，发现优化器最终选择的执行计划也不同。

统计信息。对于ANALYZE TABLE的使用，参见第七章对于表与索引的详细介绍，这里只给出ANALYZE TABLE使用的例子（在以后读完第七章后再回过头来看）：

## 2.Fast Index Creation

MySQL 5.5版本之前（不包括5.5）存在的一个普遍被人诟病的问题是MySQL数据库对于索引的添加或者删除的这类DDL操作，MySQL数据库的操作过程为：

☐首先创建一张新的临时表，表结构为通过命令ALTER TABLE新定义的结构。

☐然后把原表中数据导入到临时表。

☐接着删除原表。

☐最后把临时表重名为原来的表名。

可以发现，若用户对于一张大表进行索引的添加和删除操作，那么这会需要很长的时间。更关键的是，若有大量事务需要访问正在被修改的表，这意味着数据库服务不可用。而这对于Microsoft SQL Server、Oracle等数据库DBA来说，MySQL数据库的index维护就显得非常"不能接受"。

InnoDB存储引擎从InnoDB 1.0.x版本开始支持一种称为Fast Index Creation（快速索引创建）的索引创建方式——简称FIC。

对于辅助索引的创建，InnoDB存储引擎会对创建索引的表加上一个S锁。在创建的过程中，不需要重建表，因此速度较之前提高很多，并且数据库的可用性也得到了提高。删除辅助索引操作就更简单了，InnoDB存储引擎只需更新内部视图，并将辅助索引的空间标记为可用，同时删除MySQL数据库内部视图上对该表的索引定义即可。

这里需要特别注意的是，临时表的创建路径是通过参数tmpdir进行设置的。用户必须保证tmpdir有足够的空间可以存放临时表，否则会导致创建索引失败。

由于FIC在索引的创建的过程中对表加上了S锁，因此在创建的过程中只能对该表进行读操作，若有大量的事务需要对目标表进行写操作，那么数据库的服务同样不可用。此外，FIC方式只限定于辅助索引，对于主键的创建和删除同样需要重建一张表。

## 3.Online Schema Change

Online Schema Change（在线模式变更，简称OSC）是首先由Facebook工程师实现的DDL变更方式。由于当时Facebook的MySQL版本还未实现"秒级"加字段，而用工具的方式来规避通常锁表时间较长的MySQL原生模式DDL操作，是个好方案。

Facebook采用PHP语言实现了OSC，整体架构基于InnoDB存储引擎。自此，各种OSC工具（Facebook工程师Vamsi Ponnekanti）首先实现了OSC的一种算法，叫作The openarkkit toolkit oak-online-alter-table。整个OSC分为以下步骤。

□init：即初始化阶段，检查原表是否满足操作条件，比如是否存在主键、是否存在触发器等。

□createCopyTable：创建和原表结构一样的新表。

□alterCopyTable：对创建的新表执行ALTER TABLE，进行相应的模式变更。

□createDeltasTable：创建deltas表，作用是在数据记录操作过程中，将变更过程中的DML操作记录到该createDeltasTable中。

□createTriggers：创建触发器，INSERT、UPDATE、DELETE对应的三种类型的触发器。触发器的目的就是记录deltas表。

□startSnpshotXact：开始OSC操作的事务。

□selectTableIntoOutfile：将原表数据写出到文件。这里的实现方式比较有意思，通过分片（chunked）方式把原表数据写出到多个文件，然后加载到copy表。分片大小可以定制，默认值是500 000。

□dropNCIndexs：在数据导入到新表前，删除新表中的所有非唯一索引。

□loadCopyTable：把写出的文件数据导入到新表。

□replayChanges：在OSC操作过程中，对原表的DML操作产生的数据变更不会丢失，将这些数据变更应用到新表。这些数据就记录在deltas表中。

□recreateNCIndexes：重新创建非唯一索引。

□replayChanges：再次将原表的DML操作产生的数据变更应用到新表（配合触发器），这是为了确保数据一致。

□swapTables这个阶段出现问题的几率不是很高，此时2张表的结构已经保持同步了，只是在进行表的重命名操作。

在使用的过程中，OSC也暴露出一些问题，比如原始OSC（PHP版本）在对约2200万数据的MySQL InnoDB表进行操作时，DBA发现它产生了大量的事务操作，这使得InnoDB日志增长过快。

此外，OSC在早期的PHP版本中，对触发器进行了操作，这就要求被操作的表不能含有任何触发器，不过在原始OSC中有一个设置SET sql_bin_log=0，在此期间将不会记录到slave，这样就造成了主库与从库的数据不一致。

## 4.Online DDL

不论FIC还是InnoDB的热备插件，都需要一个临时表来进行中间过程的存储，这都需要耗费一定的磁盘空间，而且在整个DML阶段，OSC的使用都存在一个触发器，这样就会对性能有一定的影响。MySQL 5.6版本引入了Online DDL，它允许在创建索引和修改表的过程中执行查询操作和INSERT、UPDATE、DELETE这类DML操作，这就大大提升了MySQL在生产环境的可用性。

对于索引的创建与删除，修改列的DDL，可以通过指定"操作方式"来对其操作进行调整：

□是否允许读写操作。

□临时表的使用。

□临时日志文件的使用。

□执行的方式。

下面是ALTER TABLE语句操作索引时所使用的语句格式：

```
ALTER  TABLE tbl_name

|ADD{INDEX|KEY}[index_name]

[index_type](index_col_name,...)[index_option]...

ALGORITHM[=]{DEFAULT|INPLACE|COPY}

LOCK[=]{DEFAULT|NONE|SHARED|EXCLUSIVE}
```

ALGORITHM指定了执行的方式，默认是COPY，这和MySQL 5.1版本之前是一致的，即创建一张临时表。INPLACE是尝试一些操作可以不需要创建临时表，DEFAULT表示根据参数old_alter_table来判断是INPLACE还是COPY。上面的示例中参数为OFF，故采用INPLACE的方式创建。

---

mysql＞SELECT@@version\G;

***************************1.row***************************

@@version:5.6.6-m9

1 row in set(0.00 sec)

mysql＞SHOW VARIABLES LIKE'old_alter_table'\G;

***************************1.row***************************Variable_name:old_alter_table

Value:OFF

1 row in set(0.00 sec)

---

LOCK部分为索引创建或删除时对表添加锁的情况，可有的选择为：

（1）NONE

执行索引创建或者删除操作时，对目标表不添加任何的锁，即事务仍然可以进行读写操作，不会收到阻塞。因此这种模式可以获得最大的并发度。

（2）SHARE

这和之前的FIC类似，对目标表添加共享锁，但是允许事务进行读操作，但是会阻塞事务的写操作。若存储引擎不支持SHARE模式，会返回一个错误信息。

（3）EXCLUSIVE

在EXCLUSIVE模式下，对目标表添加排它锁，即事务中读写操作都被阻塞，因此会阻塞所有的事务，这和COPY对表上锁的模式类似。如果执行COPY这种创建索引的方式，

（4）DEFAULT

DEFAULT：由当前有效的系统变量决定；NONE：如果可以使用就对表加锁；SHARE：如果支持在线就使用它；EXCLUSIVE：强制使用DEFAULT，或者对于不支持在线操作的DDL采用排他锁。

InnoDB存储引擎的Online DDL在执行的过程中允许并发执行对表的INSERT、UPDATE、DELETE等DML操作。创建索引的过程中产生的对表的记录修改，会记录到innodb_online_alter_log_max_size参数控制的日志中去（默认128MB）。如果在创建索引的过程中修改的数据量超过了参数innodb_online_alter_log_max_size的限制，将会抛出类似如下的错误。

---

Error:1799SQLSTATE:HY000(ER_INNODB_ONLINE_LOG_TOO_BIG)

Message:Creating index'idx_aaa'required more than'innodb_online_alter_log_max_size'bytes of modification log.Please try again.

---

要处理这种情况，可以尝试调大innodb_online_alter_log_max_size参数值来允许更多的表记录被修改，或者在执行ALTER TABLE时指定SHARE模式来使用共享表锁，此时在创建索引的过程中不允许对表进行DML更改操作。

本节介绍了什么是在线的Online DDL，以及如何使用在线工具，在下面的实例中将实际演示在线变更SQL语句对数据库和应用造成的影响。

# 5.5 Cardinality值

## 5.5.1 什么是Cardinality

并不是在所有的查询条件中出现的列都需要添加索引。对于什么时候添加B+树索引，一般的经验是，在访问表中很少一部分时使用B+树索引才有意义。对于性别字段、地区字段、类型字段，它们可取值的范围很小，称为低选择性。如

```
SELECT*FROM student WHERE sex='M'
```

按性别进行查询时，可取值的范围一般只有'M'、'F'。因此上述SQL语句得到的结果可能是各占50%的行。这时添加B+树索引是完全没有必要的。相反，如果某个字段的取值范围很广，几乎没有重复，即属于高选择性，则此时使用B+树索引是最适合的。例如，对于姓名字段，基本上在一个应用中不允许重名的出现。

怎样查看索引是否是高选择性的呢？可以通过SHOW INDEX结果中的列Cardinality来观察。Cardinality值非常关键，表示索引中不重复记录数量的预估值。同时需要注意的是，Cardinality是一个预估值，而不是一个准确值，基本上用户也不可能得到一个准确的值。在实际应用中，Cardinality/n_rows_in_table应尽可能地接近1。如果非常小，那么用户需要考虑是否还有必要创建这个索引。故在访问高选择性属性的字段并从表中取出很少一部分的行时，对这个字段添加B+树索引是非常有必要的。如

```
SELECT*FROM member WHERE usernick='David'
```

表member大约有500万行数据。usernick字段表示用户昵称，基本上昵称是唯一的。David这个昵称的用户只有一个。因此

```
mysql）EXPLAIN SELECT*FROM member
-）WHERE usernick='David'\G;
***************************1.row***************************
id:1
```

select_type:SIMPLE

table:member

type:const

possible_keys:usernick

key:usernick

key_len:62

ref:const

rows:1

Extra:

1 row in set(0.00 sec)

---

可以看到，通过usernick字段的索引就找到了需要的数据，该SQL语句的执行性能是非常高的。

# 5.5.2 InnoDB存储引擎的Cardinality统计

索引中不重复记录数量Cardinality的预估值，但是这个Cardinality的值需要进行统计的，更新。但是需要注意的是，数据库对于Cardinality的统计都是通过采样（MySQL不同版本实现不同）的方法来完成的，B+树索引的更新又不是实时的Cardinality值进行更新的，因此此值并不是实时更新的，即非实时更新。

而在生产环境中，索引的更新操作可能是非常频繁的。如果每次发生变更时就去更新Cardinality值，那么将会给数据库带来很大的负担。另外需要考虑的一个因素是，如果表50G，那么需要对这个表的Cardinality进行统计，这是需要扫描表中所有数据。因此，数据库对于Cardinality的统计都是通过采样（Sample）的方法来完成的。

在InnoDB存储引擎中，Cardinality统计信息的更新发生在两个操作中：INSERT和UPDATE。根据前面的叙述，不可能在每次发生INSERT和UPDATE时就去更新Cardinality信息，这样会增加数据库系统的负荷，同时对于大表的统计，需要时间上的开销。因此，InnoDB存储引擎内部对更新Cardinality信息的策略为：

□表中1/16的数据已发生过变化。

□stat_modified_counter＞2 000 000 000。

第一种策略为，若表中的Cardinality值。这里的1/16是因为采样的数据量为16个叶子结点。如果每次Cardinality值都更新，那么将会对数据库系统造成负荷，同时会影响系统运行性能。而第二种情况，考虑的是，InnoDB存储引擎内部有一个计数器stat_modified_counter，用来表示发生变化的次数，当stat_modified_counter大于2 000 000 000时，则同样需要更新Cardinality信息。

接着来看InnoDB存储引擎内部是怎样来进行Cardinality信息的统计和更新操作的呢？在InnoDB存储引擎中，默认8个叶子结点（Leaf　Page）进行采样。采样的过程如下：

□取得B+树索引中叶子结点的数量，记为A。

□随机取得B+树索引中8个叶子结点，统计每个页不同记录的个数，即为P1、P2、…、P8。

□的平均值作为最后的Cardinality，即：Cardinality=（P1+P2+…+P8）*A/8。

从上述的公式中可以发现，在InnoDB存储引擎中，Cardinality统计信息的更新发生在两个操作中：INSERT和UPDATE。但是，并不是每次发生时都会更新Cardinality统计信息，这样会增加负担，而且对于Cardinality信息是一个精确值也没有太大的意义。

```
SHOW INDEX FROM OrderDetails
```

因此，通过SQL命令得到的MySQL数据库中Cardinality值是一个大概的值，如下面的图5-20所示。



图 5-20  执行语句「SHOW INDEX FROM OrderDetails」的结果

因此，如果表的数据进行了大量的INSERT、UPDATE操作后，对表OrderDetails进行观察，再次执行语句SHOW INDEX FROM后，会发现Cardinality值发生了变化，如下面的图5-21所示。

| Table | Non_unique | Key_name | Seq_in_index | Column_name | Collation | Cardinality | Sub_part | Packed | Null | Index_type | Comment | Index_comment |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| orderdetails | 0 | PRIMARY | 1 | OrderID | A | 2192 | NULL | NULL | | BTREE | | |
| orderdetails | 0 | PRIMARY | 2 | ProductID | A | 2192 | NULL | NULL | | BTREE | | |
| orderdetails | 1 | OrderID | 1 | OrderID | A | 2192 | NULL | NULL | | BTREE | | |
| orderdetails | 1 | OrdersOrder_Details | 1 | OrderID | A | 2192 | NULL | NULL | | BTREE | | |
| orderdetails | 1 | ProductID | 1 | ProductID | A | 168 | NULL | NULL | | BTREE | | |
| orderdetails | 1 | ProductsOrder_Details | 1 | ProductID | A | 168 | NULL | NULL | | BTREE | | |

图 5-21 　　　执行命令SHOW INDEX FROM OrderDetails的结果

通过上图可以发现，在运行SHOW INDEX FROM命令后，表OrderDetails中列Cardinality的值还是保持没有变化，即表OrderDetails的统计信息并没有更新。由于每次对Cardinality值的统计都是通过8个叶子节点来完成，故统计的叶子节点可能不同，所以对于Cardinality值可能会有不同的值。因此这是InnoDB存储引擎Bug吗？其实这并不是一个Bug。

更新索引统计信息的方法有两种，可能会导致Cardinality值即使通过命令进行设置也不会发生变化。例如每次通过8个叶子节点来进行采样，但是采样的过程中，数据页都没有发生过读取、写入操作，即Cardinality值的更新是有条件的。

当InnoDB 1.2版本时，可以通过参数innodb_stats_sample_pages用来设置统计Cardinality时每次采样页的数量，默认值为8。同时，参数innodb_stats_method用来判断如何对待索引中出现的NULL值记录。该参数的默认值为nulls_equal，表示将NULL值记录视为相等的记录。其有效值还有nulls_unequal、nulls_ignored，分别表示将NULL值记录视为不相等的记录和忽略NULL值记录。例如某页中索引记录为NULL、NULL、1、2、2、3、3、3，在参数innodb_stats_method的默认设置下，表中该列的Cardinality为4；若参数innodb_stats_method为nulls_unequal，则该列的Caridinality为5；若参数innodb_stats_method为nulls_ignored，则Cardinality为3。

此外，在SQL语句ANALYZE TABLE、SHOW TABLE STATUS、SHOW INDEX以及访问INFORMATION_SCHEMA架构下的表TABLES和STATISTICS时，都会导致InnoDB存储引擎去重新计算索引的Cardinality值。若表中的数据量非常大，并且表中存在多个辅助索引时，执行上述这些操作可能会非常慢。虽然用户可能并不希望去更新Cardinality值。

InnoDB1.2版本之前，使用采样的方法来计算Cardinality。采样的过程如下（由图5-3所示）：

表 5-3　InnoDB 1.2 新增参数

| 参数 | 说明 |
|---|---|
| innodb_stats_persistent | 是否将命令 ANALYZE TABLE 计算得到的 Cardinality 值存放到磁盘上。若是，则这样做的好处是可以减少重新计算每个索引的 Cardinality 值，例如当 MySQL 数据库重启时。此外，用户也可以通过命令 CREATE TABLE 和 ALTER TABLE 的选项 STATS_PERSISTENT 来对每张表进行控制。<br><br>默认值：OFF |
| innodb_stats_on_metadata | 当通过命令 SHOW TABLE STATUS、SHOW INDEX 及访问 INFORMATION_SCHEMA 架构下的表 TABLES 和 STATISTICS 时，是否需要重新计算索引的 Cardinality 值。<br><br>默认值：OFF |
| innodb_stats_persistent_sample_pages | 若参数 innodb_stats_persistent 设置为 ON，该参数表示 ANALYZE TABLE 更新 Cardinality 值时每次采样页的数量。<br><br>默认值：20 |
| innodb_stats_transient_sample_pages | 该参数用来取代之前版本的参数 innodb_stats_sample_pages，表示每次采样页的数量。<br><br>默认值为：8 |

# 5.6 B+□□□□□□

## 5.6.1 □□□□□□B+□□□□□□□

□□□□□B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Think Different□

□□□□1□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLTP□OLAP□□□□□OLTP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□10□□□□□□□□□□□□□□□□□□□□1□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLTP□□□□□□□□□□□□□□□□□□□B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□B+□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Hash Join□□□□□□□□□□□□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□□OLAP□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 5.6.2 联合索引

联合索引是指对表上的多个列进行索引。前面讨论的情况都是只对表上的一个列进行索引。联合索引的创建方法与单个索引创建的方法一样，不同之处仅在于有多个索引列。例如，下面代码创建了一张t表，并且索引idx_a_b是联合索引，联合的列为（a，b）。

---

CREATE TABLE t(

a INT,

b INT,

PRIMARY KEY(a),

KEY idx_a_b(a,b)

)ENGINE=INNODB

---

从本质上来说，联合索引也是一棵B+树，不同的是联合索引的键值的数量不是1，而是大于等于2。接着来讨论两个整型列组成的联合索引，假定两个键值的名称分别为a、b，如图5-22所示。



图 5-22 多个键值的B+树

图5-22中显示了(a,b)的B+树索引，可以进行(a,b)的联合查询。需要注意的是，对于联合索引(a,b)，B+树索引是根据联合索引的值进行排序的，因此(1，1)、(1，2)、(2，1)、(2，4)、(3，1)、(3，2)是按照a、b的顺序进行存放的。

因此，对于查询SELECT*FROM TABLE WHERE a=xxx and b=xxx，显然是可以使用(a,b)这个联合索引的。对于单个的a列查询SELECT*FROM TABLE WHERE a=xxx，也可以使用这个(a,b)索引。但对于b列的查询SELECT*FROM TABLE WHERE b=xxx，则不可以使用这棵B+树索引。可以发现叶子节点上的b值为1、2、1、4、1、2，显然不是排序的，因此对于b列的查询使用不到(a,b)索引。

联合索引的第二个好处是已经对第二个键值进行了排序处理。例如，在很多情况下应用程序都需要查询某个用户的购物情况，并按照时间进行排序，最后取出最近三次的购买记录，这时使用联合索引可以避免多一次的排序操作，因为索引本身在叶子节点已经排序了。来看一个例子，首先根据如下代码来创建测试表buy_log：

CREATE TABLE buy_log(

userid INT UNSIGNED NOT NULL,

buy_date DATE

)ENGINE=InnoDB;

INSERT INTO buy_log VALUES(1,'2009-01-01');

INSERT INTO buy_log VALUES(2,'2009-01-01');

INSERT INTO buy_log VALUES(3,'2009-01-01');

INSERT INTO buy_log VALUES(1,'2009-02-01');

INSERT INTO buy_log VALUES(3,'2009-02-01');

INSERT INTO buy_log VALUES(1,'2009-03-01');

INSERT INTO buy_log VALUES(1,'2009-04-01');

ALTER TABLE buy_log ADD KEY(userid);

ALTER TABLE buy_log ADD KEY(userid,buy_date);

上述代码建立了两个索引来进行比较。两个索引都包含了userid字段。如果只对于userid进行查询，如：

SELECT*FROM buy_log WHERE userid=2;

查询结果执行计划如图5-23所示。



图 5-23 仅根据索引userid进行的查询

在图5-23中可以看到，possible_keys这一列显示了查询中有两个索引可选，即userid和联合索引（userid，buy_date）。但优化器最终选择使用的是索引userid，因为该索引的叶节点包含单个键值，所以理论上一个页能存放的记录应该更多。

接着假设要取出userid为1的用户最近3次的购买记录，用SQL语句来表示如下，其执行计划如图5-24所示。

---

SELECT*FROM buy_log

WHERE userid=1 ORDER BY buy_date DESC LIMIT 3

---



图 5-24 SQL语句的执行计划

在这种情况下，上述SQL语句可以利用userid，buy_date联合索引（userid，buy_date）来进行优化。因为userid，buy_date联合索引即userid_2已经对buy_date进行了排序，所以无需再对buy_date进行额外的排序操作。通过EXPLAIN命令，用户可以观察到优化器实际上选择了userid，而这次显示了。如图5-25所示。



图 5-25 优化器userid索引进行查询

若Extra列里还出现了Using filesort，那就说明存在隐患了。比如下面这个例子中是先根据buy_date排序，再根据userid、buy_date联合排序。

我们再来看一个例子，假设你有a、b两个索引，但在a、b两个字段上排序，排序时无法使用到索引：

---

SELECT...FROM TABLE WHERE a=xxx ORDER BY b

---

接着假设你有a、b、c三个索引，下面这两种情况下都可以用上索引进行排序：

---

SELECT...FROM TABLE WHERE a=xxx ORDER BY b

SELECT...FROM TABLE WHERE a=xxx AND b=xxx ORDER BY c

---

但是下面这种情况下无法用上索引进行排序，因为它违反了最左匹配原则，会出现filesort，因为在索引里，是先按a、c排序的。

---

SELECT...FROM TABLE WHERE a=xxx ORDER BY c

---

# 5.6.3 □□□□

InnoDB□□□□□□□□□□□□covering index□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□IO□□□□

□□ □□□□□□□□□□□InnoDB Plugin□□□□□□□□□□□□□□□□InnoDB□□□□□1.0 □□□□□MySQL□□□□□□□5.0□□□□□□InnoDB□□□□□□□□□□□□□□□□

□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□primary key1□primary key2□...□key1□key2□...□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
SELECT key2 FROM table WHERE key1=xxx□

SELECT primary key2,key2 FROM table WHERE key1=xxx□

SELECT primary key1,key2 FROM table WHERE key1=xxx□

SELECT primary key1,primary key2□key2 FROM table WHERE key1=xxx□
```

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□buy_log□□□□□□□□□□□□□□

```
SELECT COUNT(*)FROM buy_log;
```

InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□buy_log□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□IO□□□□□□□□□□□□□□□5-26□□□□

| id | select_type | table | type | possible_keys | key | key_len | ref | rows | Extra |
|----|-------------|-------|------|---------------|-----|---------|-----|------|-------|
| 1 | SIMPLE | buy_log | index | NULL | userid | 4 | NULL | 7 | Using index |

□ 5-26 COUNT□*□□□□□□□□□□

从图5-26中可以看到，possible_keys为NULL，但是实际执行时用到了userid索引，同时在Extra中显Using index，表示索引覆盖了查询。接着来看如下语句：

如果只是统计某个用户购买的情况，如a、b所示。这也符合联合索引的使用方式，即b字段的查询条件对应于联合索引最左侧字段的查询条件，这样仍然能够使用这个联合索引：

---

```
SELECT COUNT(*)FROM buy_log

WHERE buy_date>='2011-01-01'AND buy_date<'2011-02-01'
```

---

对buy_log上（userid，buy_date）的联合索引，此时仅对于b列进行条件的查询，一般情况下是不能进行该SQL语句的优化的，因为在联合索引中列顺序问题，但是因为是统计操作，并且可以进行索引覆盖，因此优化器会选择该联合索引，如图5-27所示。



图 5-27   索引覆盖用于优化统计操作

从图5-27中可以看到，possible_keys依然为NULL，但是key为userid_2，即表示（userid，buy_date）的联合索引，同时Extra列的信息为Using index，表示进行了索引覆盖操作。

# 5.6.4 索引的使用及其注意事项

索引使用最常用的是使用EXPLAIN来跟踪SQL语句是否使用了索引及索引的使用是否合适。下面展示一个实际应用例子。这个例子所采用的表是一个两个表进行JOIN后经过细粒度化产生的。

---

SELECT*FROM orderdetails

WHERE orderid□10000 and orderid□102000;

---

上面这条SQL语句意欲取出订单号10000到某范围的订单信息。通过SHOW INDEX FROM orderdetails命令查看索引,如图5-28所示。



图 5-28 表orderdetails的索引信息

从图中可见orderdetails表在OrderID、ProductID上有联合索引,但在单独的OrderID列上也有索引。对于上面的SQL语句,查询优化器利用OrderID上这个单独的索引来完成查询。通过EXPLAIN命令查看的结果,可以发现在OrderID上使用了索引,如图5-29所示。

图 5-29　执行查询语句SQL的结果

从possible_keys可以看出，优化器有三个索引PRIMARY、OrderID、OrdersOrder_Details可供选择，而优化器最终选择了主键索引PRIMARY。出现这个问题的原因是，table scan的开销要小于index scan的。

在上面这个示例中，由于返回的记录数较多，如果通过OrderID辅助索引来进行访问，访问的数据量太大，通过OrderID辅助索引来访问数据所产生的逻辑读，要远远大于通过OrderID索引来访问数据所产生的开销。虽然OrderID辅助索引中数据顺序存储，但是再根据主键聚集索引去查找数据是随机的。因此成为了随机的读操作。故在索引返回的记录数超过表中数据20%的情况下，优化器通常会选择通过聚集索引来查找数据。

当然，并不是所有情况下通过辅助索引查询的效率都要低于聚集索引。在有些情况下，通过辅助索引查询的效率是很高的，但是优化器有时会将这部分情况也错误地判断为聚集索引查询效率更高。这时，就可以强制通过FORCE INDEX来进行索引的选择了。例如：

```
SELECT*FROM orderdetails FORCE INDEX(OrderID)

WHERE orderid□10000 and orderid□102000;
```

执行结果如图（图5-30所示）



图 5-30　强制使用索引查询

# 5.6.5 □□□□

MySQL□□□□□□□□□□□INDEX HINT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□INDEX HINT□

□MySQL□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□
MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
DBA□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□

□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□
□□□□□Range□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□□□□Index
Hint□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□MySQL□□□□Index Hint□□□□□□□

```
tbl_name[[AS]alias][index_hint_list]

index_hint_list:

index_hint[,index_hint]...

index_hint:

USE{INDEX|KEY}

[{FOR{JOIN|ORDER BY|GROUP BY}]([index_list])

|IGNORE{INDEX|KEY}

[{FOR{JOIN|ORDER BY|GROUP BY}](index_list)

|FORCE{INDEX|KEY}

[{FOR{JOIN|ORDER BY|GROUP BY}](index_list)

index_list:

index_name[,index_name]...
```

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□t□□□□□□□□□□□□

```
CREATE TABLE t(

a INT,
```

```
b INT,

KEY(a),

KEY(b)

)ENGINE=INNODB;

INSERT INTO t SELECT 1,1;

INSERT INTO t SELECT 1,2;

INSERT INTO t SELECT 2,3;

INSERT INTO t SELECT 2,4;

INSERT INTO t SELECT 1,2;
```

## 接着来看如下的SQL语句：

```
SELECT*FROM t WHERE a=1 AND b=2;
```

## 通过EXPLAIN得到如图5-31所示的执行计划。

| id | select_type | table | type | possible_keys | key | key_len | ref | rows | Extra |
|----|-------------|-------|------|---------------|-----|---------|-----|------|-------|
| 1 | SIMPLE | t | index_merge | a,b | b,a | 5,5 | NULL | 1 | Using intersect(b,a); Using where; Using index |

## 图 5-31 SQL语句的执行计划

从5-31中的possible_keys可以看到，此时SQL优化器可以选择表t上的a、b索引，最终选择的索引为key所示的a、b索引，即MySQL数据库将a、b两个索引进行索引合并，在Extra列中可看到Using intersect（b、a）的提示信息，表示根据两个索引得到的结果进行求交的数学运算。

### 用户可以通过USE INDEX提示，强制使用a这一个索引，如：

```
SELECT*FROM t USE INDEX(a)WHERE a=1 AND b=2;
```

## 这时的执行计划如图5-32所示。

| id | select_type | table | type | possible_keys | key | key_len | ref | rows | Extra |
|---|---|---|---|---|---|---|---|---|---|
| 1 | SIMPLE | t | ALL | a | NULL | NULL | NULL | 5 | Using where |

图 5-32 使用USE INDEX未使用索引

从图中可以看到，虽然看到可以使用索引a，但是优化器最终选择通过扫描聚集索引，即直接进行全表扫描来得到数据。因此强制使用索引可以使用FORCE INDEX，如下语句所示：

```
SELECT*FROM t FORCE INDEX(a)WHERE a=1 AND b=2;
```

上述语句的执行计划如图5-33所示。

| id | select_type | table | type | possible_keys | key | key_len | ref | rows | Extra |
|---|---|---|---|---|---|---|---|---|---|
| 1 | SIMPLE | t | ref | a | a | 5 | const | 3 | Using where |

图 5-33 使用FORCE INDEX强制使用索引

这时优化器的最终选择和用户指定的索引是一致的。因此，如果用户确定需要使用某个索引来得到数据库，可以使用FORCE INDEX，而不是USE INDEX。

# 5.6.6 Multi-Range Read优化

MySQL5.6版本开始支持Multi-Range Read（MRR）优化。Multi-Range Read优化的目的就是为了减少磁盘的随机访问，并且将随机访问转化为较为顺序的数据访问，这对于IO-bound类型的SQL查询语句可带来性能极大的提升。Multi-Range Read优化可适用于range，ref，eq_ref类型的查询。

MRR优化有以下几个好处：

❑MRR使数据访问变得较为顺序。在查询辅助索引时，首先根据得到的查询结果，按照主键进行排序，并按照主键排序的顺序进行书签查找。

❑减少缓冲池中页被替换的次数。

❑批量处理对键值的查询操作。

对于InnoDB和MyISAM存储引擎的范围查询和JOIN查询操作，MRR的工作方式如下：

❑将查询得到的辅助索引键值存放于一个缓存中，这时缓存中的数据是根据辅助索引键值排序的。

❑将缓存中的键值根据RowID进行排序。

❑根据RowID的排序顺序来访问实际的数据文件。

此外，若InnoDB存储引擎或者MyISAM存储引擎的缓冲池不是足够大，即不能存放下一张表中的所有数据，此时频繁的离散读操作还会导致缓存中的页被替换出缓冲池，然后又不断地被读入缓冲池。若是按照主键顺序进行访问，则可以将此重复行为降为最低。如下面这句SQL语句：

---

```
SELECT*FROM salaries WHERE salary＞10000 AND salary＜40000;
```

---

salary上有一个辅助索引idx_s，因此除了通过辅助索引查找键值外，还需要通过书签查找来进行对整个表数据的查询。当不启用Multi-Range Read特性时，看到的执行计划如图5-34所示。

图 5-34　未开启Multi-Range Read时执行计划

未开启Mulit-Range Read的执行计划，其列Extra中为Using index condition，而非我们预期的Using MRR，如图5-35所示。



图 5-35　开启Multi-Range Read时执行计划

对于上述两句相同的SQL语句，其执行时间的对比如表5-4所示。

表5-4　是否启用 Multi-Range Read 的执行时间对比

| | 执行时间 ( 秒 ) |
| --- | --- |
| 不使用 Multi-Range Read | 43.213 |
| 使用 Multi-Range Read | 4.212 |

可以看到，启用该特性后，执行时间大约快了10倍。可见，Multi-Range Read特性对于磁盘的离散访问优化非常大。

**③　索引下推**　在MySQL数据库中还支持一种称为Index Condition Pushdown的优化方式。同样，该优化支持range、eq_ref、

启用Multi-Range Read特性后，优化器会根据查询条件，将随机的访问转化为较为顺序的数据访问请求，从而大幅提高访问数据的效率。如：

---

```
SELECT*FROM t

WHERE key_part1＞=1000 AND key_part1＜2000

AND key_part2=10000;
```

---

表t有（key_part1，key_part2）的联合索引，因此索引根据key_part1、key_part2的位置关系进行排序。若没有Multi-Read Range，此时查询类型为Range，SQL优化器会先将key_part1大于1000且小于2000的数据都取出，即使key_part2不等于1000。待取出行数据后再根据key_part2的条件进行过滤。这会导致无用数据被取出。如果有大量数据且key_part2不等于1000，则启用Mulit-Range Read特性会使性能有巨大的提升。

倘若启用Multi-Range Read特性，优化器会先将查询条件进行拆分，然后再进行数据查询。就上述查询语句而言，优化器会将查询条件拆分为（1000，1000），（1001，1000），（1002，1000），...，（1999，1000）。最后再根据这些拆分出的条件进行数据的查询。

下面是另一个使用的例子，如：

---

```
SELECT*FROM salaries

WHERE(from_date between'1986-01-01'AND'1995-01-01')

AND(salary between 38000 and 40000);
```

---

若启用了Multi-Range Read特性，则执行计划如图5-36所示。

| id | select_type | table | type | possible_keys | key | key_len | ref | rows | Extra |
|---|---|---|---|---|---|---|---|---|---|
| 1 | SIMPLE | salaries | range | idx_s | idx_s | 4 | NULL | 210740 | Using index condition; Using MRR |

图 5-36 启用Multi-Range Read后的执行计划

对salaries表根据salary列构建idx_s索引，执行下面SQL语句，优化器采用Multi-Range Read优化机制，可以从下面执行计划的Extra列中看到Using MRR的信息。

是否采用Multi-Range Read优化机制，可由optimizer_switch中的两个flag进行控制，当mrr为on时，表示启用Multi-Range Read，当mrr_cost_based为时，表示采用cost based的方式进行选择是否mrr，如果mrr设为on，mrr_cost_based设为off，则表示始终开启Multi-Range Read优化机制，如下所示为强制开启Multi-Range Read优化机制的命令示例。

---

```
mysql〉SET@@optimizer_switch='mrr=on,mrr_cost_based=off';

Query OK,0 rows affected(0.00 sec)
```

---

通过read_rnd_buffer_size来设置缓存的大小，当缓存不够时，则执行器会将已经缓存的RowID进行排序，并通过排好序的RowID去缓存数据，其默认值为256K。

---

```
mysql〉SELECT@@read_rnd_buffer_size\G;

***************************1.row***************************

@@read_rnd_buffer_size:262144

1 row in set(0.00 sec)
```

---

# 5.6.7　Index Condition Pushdown（ICP）优化

和Multi-Range Read优化和Index Condition Pushdown优化是MySQL 5.6开始支持的一种根据索引进行查询的优化方式。之前MySQL数据库版本不支持Index Condition Pushdown时，进行索引查询时，首先根据索引来查找记录，然后再根据WHERE条件来过滤记录。在支持Index Condition Pushdown后，MySQL数据库会在取出索引的同时，判断是否可以进行WHERE条件的过滤，也就是将WHERE的部分过滤操作放在了存储引擎层。在某些查询下，可以大大减少上层SQL层对记录的索引（fetch），从而提高数据库的整体性能。

Index Condition Pushdown优化支持range、ref、eq_ref、ref_or_null类型的查询，当前支持MyISAM、InnoDB存储引擎。若支持Index Condition Pushdown优化，则会在执行计划的列Extra看到Using index condition提示。

此外　NDB Cluster存储引擎有Engine Condition Pushdown，千万不要混淆，这两者是不同的。"Index"的Condition Pushdown是索引下推，是将Condition Pushdown放在了存储引擎层的索引中。而早在MySQL 5.1时，对NDB Cluster存储引擎支持了Engine Condition Pushdown。

假设某张表有联合索引(zip_code，last_name，firset_name)，并且查询语句如下：

```
SELECT*FROM people

WHERE zipcode='95054'

AND lastname LIKE'%etrunia%'

AND address LIKE'%Main Street%';
```

对于前面的例子，MySQL数据库可以通过索引来定位zipcode等于95 054的记录，但是索引对于WHERE条件的lastname LIKE'%etrunia%'AND address LIKE'%Main Street%'没有任何帮助。若不支持Index Condition Pushdown优化，则数据库需要先通过索引取出所有zipcode等于95 054的记录，然后再过滤WHERE之后的两个条件。

若支持Index Condition Pushdown优化，则在索引取出时，就会进行WHERE条件的过滤，然后再去获取记录。这将大大减少WHERE条件的判断，从而提高整体的查询效率，减少上层SQL层的查询。

```
SELECT*FROM salaries

WHERE(from_date between'1986-01-01'AND'1995-01-01')

AND(salary between 38000 and 40000);
```

若不启用Multi-Range Read特性，则执行计划如图5-37所示。



图 5-37    不启用Multi-Range Read特性的执行计划

此时可以在Extra（Using index condition）看到又使用了索引idx_s，并使用了Index Condition Pushdown优化。这里索引根据列(emp_no、from_date)进行排序，因此使用idx_s也可过滤列from_date数据，而非只过滤单个主键。

表5-5显示了MySQL 5.5与MySQL 5.6中该SQL语句的执行时间，以及是否支持MRR新特性的区别。

表 5-5    MySQL 5.5 和 MySQL 5.6 中是否启用 Index Condition Pushdown 的执行时间对比

|  | 执行时间（秒） |
| --- | --- |
| MySQL 5.5 | 46.738 |
| MySQL 5.6 with ICP | 37.924 |
| MySQL 5.6 with ICP & MRR | 7.816 |

通过查询优化器的这些特性，我们可以写出更高效的查询语句。例如，Index Condition Pushdown（索引条件下推）可以使MySQL 5.5的某些查询速度提升23%，同时使用了Mulit-Range Read（多范围读取），可以使查询速度提升400%以上。

## 5.7　哈希表

字典、集合乃至向量，都支持在常数时间内，即以O（1）复杂度实现高效查找，但其对存储空间的消耗往往很大。比如即便是128GB的主存空间，也不足以支持对长整数范围内关键码的直接寻址。反之，一旦我们不再奢望如此高的查找效率，转而仅要求平均意义上的O（1）复杂度，则可以更好地控制对空间的消耗。

## 5.7.1　原理

哈希表（Hash Table）即是这样一种思想的具体实现，其基本原理如下。假设所有可能出现的关键码，构成全集U（不妨暂且假设为自然数的一个子集），且各关键码互异。哈希表就是全集的一个子集U＝{0，1，…，m-1} [1]，用它们作为下标，即可直接访问对应的数据项。

具体地，可以通过一个数组T[0..m-1]来存放这些数据项，其中的每一单元称作一个桶。由于U中的关键码（如图5-38所示）总是小于m的，故对于关键码k，我们约定将其存放于第k号桶中，即T[k]＝NULL。

图 5-38  直接寻址法

设所有可能出现的关键字组成关键字全域U，将每个关键字映射到下标为关键字本身的元素位置。为此，需设计一个一维数组（即直接地址表）T，并使U中的每个关键字都有一个与之对应的数组元素，即关键字集合K是全域U的子集。根据关键字在表T中存放位置的不同，可以设计出直接寻址表。

哈希表实际上是基于数组的一种扩展，由一个函数h（k）来建立关键字与函数值h之间的关系，k是待查找的数据，根据h计算出全域U中的数据存放在T[0..m-1]之间，如图5-39所示。

图 5-39 散列表

只有当两个关键字有相同的散列值时才会发生冲突，散列表中一个槽对应多个关键字，图5-39中就显示了这样的情况。此时我们称发生了冲突（collision），解决冲突的一个最简单的方法就是把散列到同一个槽中的所有元素都放在一个链表中，称为链接法（chaining）。

我们把一个散列表中每个单元都看作一个指针，指向一个链表，如图5-40所示。插入j时，将j插入到它的散列值所对应的链表的头部，查找和删除j也是一样的方法。若链表为空，则置j值为NULL。

图 5-40 　用链接法解决碰撞的散列表

　　在上述散列表中，由散列函数h将关键字值映射到散列表的各个位置上。因为关键字域远大于散列表的长度，所以有可能将不同的关键字映射到散列表的同一个位置上，这时就发生了碰撞。

　　最常用的散列函数是除余法，将关键字k除以m，取其余数作为散列地址。即将关键字k除以m（槽的数目）后所得余数作为散列表地址的方法。

h(k)=k mod m

[1]　对于m的选择，经验表明，

## 5.7.2　InnoDB□□□□□□□□□□□□□

InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Page□□□□□chain□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□m□□□□□□□□□2□□□□□□□□□□□□□□□□□□□□□□□innodb_buffer_pool_size□□□□□10M□□□□□640□16KB□□□□□□□□□□□□□□□□□□□□□□□□□□640×2=1280□□□□□□□□□1280□□□□□□□□□□□□1280□□□□□□□□□□□□□□1399□□□□□□□□□□□□□1399□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□InnoDB□□□□□□□□□□□□□□□space_id□□□□□□□□□□□□□□□□□□□□□□□□□□□□16KB□□□□□□□□offset□InnoDB□□□□□space_id□□20□□□□□□□□□□space_id□offset□□□□□□□K=space_id□□20+space_id+offset□□□□□□□□□□□□□□□□□□□

# 5.7.3 □□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□DBA□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
SELECT*FROM TABLE WHERE index_col='xxx'□□□□□□□□□□□□□□□
□□□□□□□□□SHOW ENGINE INNODB STATUS□□□□□□□□□□□□□□□□□□□□
□□□□□□

---

```
mysql□SHOW ENGINE INNODB STATUS\G;

***************************1.row***************************

Status:

==================================

090922 11:52:51 INNODB MONITOR OUTPUT

==================================

Per second averages calculated from the last 15 seconds

……

-------------------------------------

INSERT BUFFER AND ADAPTIVE HASH INDEX

-------------------------------------

Ibuf:size 2249,free list len 3346,seg size 5596,

374650 inserts,51897 merged recs,14300 merges

Hash table size 4980499,node heap has 1246 buffer(s)

1640.60 hash searches/s,3709.46 non-hash searches/s

……
```

---

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

```
SELECT*FROM table WHERE index_col='xxx'
```

---

通过监控如上的一些数据，例如每秒的查询次数，每秒使用自适应哈希索引搜索的状态，可以体会到非哈希索引搜索non-hash searches/s情况下，使用hash searches:non-hash searches情况可以帮助查看当前自适应哈希索引的使用情况。

值得注意的是，哈希索引只能用来搜索等值的查询，而对于其他类型的查找操作，例如范围查找，是不能使用哈希索引的，因此这里出现了InnoDB存储引擎还提供了参数来考虑是否开启自适应哈希索引，默认值为参数innodb_adaptive_hash_index，默认值为开启，可以通过参数来禁用。

# 5.8 全文检索

## 5.8.1 概述

通过前面章节的介绍已经知道，B+树索引的特点，可以通过索引字段的前缀（prefix）进行查找。例如，对于下面的查询B+树索引是支持的：

```
SELECT*FROM blog WHERE content like'xxx%'
```

上述SQL语句可以查询博客内容以xxx开头的文章，并且只要content添加了B+树索引，就能利用索引进行快速查询。然而实际这种查询不符合用户的要求，因为在更多的情况下，用户需要查询的是博客内容包含单词xxx的文章，即：

```
SELECT*FROM blog WHERE content like'%xxx%'
```

根据B+树索引的特性，上述SQL语句即便添加了B+树索引也是需要进行索引的扫描来得到结果。类似这样的需求在互联网上搜索引擎中广泛存在。用户输入若干个关键字，搜索引擎会帮助用户定位到这个关键字的文章，这些都不是B+树索引所能很好地完成的工作。

全文检索（Full-Text Search）是将存储于数据库中的整本书或整篇文章中的任意内容信息查找出来的技术。它可以根据需要获得全文中有关章、节、段、句、词等信息，也可以进行各种统计和分析。

在之前的MySQL数据库中，InnoDB存储引擎并不支持全文检索技术。大多数的用户转向MyISAM存储引擎，这可能需要进行表的拆分，并将需要进行全文检索的数据存储为MyISAM表。这样的架构difficult带来了维护的复杂性。然而从InnoDB 1.2.x版本开始，InnoDB存储引擎开始支持全文检索，其支持MyISAM存储引擎的全部功能，并且还支持其他的一些特性。这些内容将在后面的小节中进行介绍。

从InnoDB 1.2.x版本开始，InnoDB存储引擎开始支持全文检索，其支持MyISAM存储引擎的全部功能，并且还支持其他的一些特性。这些内容将在后面的小节中进行介绍。

## 5.8.2　倒排索引

全文检索通常使用倒排索引（inverted index）来实现。倒排索引同B+树索引一样，也是一种索引结构。它在辅助表（auxiliary table）中存储了单词与单词自身在一个或多个文档中所在位置之间的映射。这通常利用关联数组实现，其拥有两种表现形式：

□inverted file index，其表现形式为{单词，单词所在文档的ID}

□full invertedindex，其表现形式为{单词，(单词所在文档的ID，在具体文档中的位置)}

以表5-6所示的全文检索表t为例来进行说明，表5-6由两列组成。

表 5-6　全文检索表 t

| DocumentId | Text | DocumentId | Text |
|---|---|---|---|
| 1 | Pease porridge hot, pease porridge cold | 4 | Some like it hot, some like it cold |
| 2 | Pease porridge in the pot | 5 | Some like it in the pot |
| 3 | Nine days old | 6 | Nine days old |

DocumentId表示进行全文检索文档的Id，Text表示存储的内容，用户需要对存储的这些文档内容进行全文查询。比如查询出现过Some这个单词的文档Id，又或者查询单个文档中出现过Some这个单词的Id和位置。

对于inverted file index的关联数组，其存储的内容如表5-7所示。

表 5-7  inverted file index 的关联数组

| Number | Text | Documents | Number | Text | Documents |
|---|---|---|---|---|---|
| 1 | code | 1, 4 | 8 | old | 3, 6 |
| 2 | days | 3, 6 | 9 | pease | 1, 2 |
| 3 | hot | 1, 4 | 10 | porridge | 1, 2 |
| 4 | in | 2, 5 | 11 | pot | 2, 5 |
| 5 | it | 4, 5 | 12 | some | 4, 5 |
| 6 | like | 4, 5 | 13 | the | 2, 5 |
| 7 | nine | 3, 6 | | | |

例如，单词"code"出现在文档1和4中，单词"days"出现在文档3和6中。我们可以看到，上表中的Documents列实际上就是一个集合。如果inverted file index存储的是文档Id，那full inverted index存储的则是(pair)，即(DocumentId，Position)，其对应的关联数组见表5-8所示。

## 表 5-8 full inverted index 的关联数组

| Number | Text | Documents | Number | Text | Documents |
|--------|------|-----------|--------|------|-----------|
| 1 | code | （1:6），（4:8） | 8 | old | （3:3），（6:3） |
| 2 | days | （3:2），（6:2） | 9 | pease | （1:1,4），（2:1） |
| 3 | hot | （1:3），（4:4） | 10 | porridge | （1:2,5），（2:2） |
| 4 | in | （2:3），（5:4） | 11 | pot | （2:5），（5:6） |
| 5 | it | （4: 3,7），（5:3） | 12 | some | （4:1,5），（5:1） |
| 6 | like | （4:2,6），（5:2） | 13 | the | （2:4），（5:5） |
| 7 | nine | （3:1），（6:1） | | | |

full inverted index□□□□□□□□□□□□□□code□□□□□□□□□1:6□□□□□□1
□□6□□□□□code□□□□□□□full inverted index□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□

# 5.8.3　InnoDB□□□□

InnoDB□□□□□□1.2.x□□□□□□□□□□□□□□□□□□full inverted index□□□□□□□InnoDB□□□□□□□□(DocumentId□Position)□□□□□"ilist"□□□□□□□□□□□□□□□□□□□□□□word□□□□□□□□□ilist□□□□□□□word□□□□□□□□□□□□□□□InnoDB□□□□□ilist□□□□□□□□Position□□□□□□□□□□Proximity Search□□MyISAM□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□word□□□□□□□□□□□□□□□Auxiliary Table□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□6□Auxiliary Table□□□□□□□□□word□Latin□□□□□□□□□

Auxiliary Table□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□FTS Index Cache□□□□□□□□□□□□□□□□□□□□□□□□□□□

FTS Index Cache□□□□□□□□□□□□word□ilist□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□FTS Index Cache□□□Auxiliary Table□□□□□□□□□InnoDB□□□□□□□□□Auxiliary Table□□□□□□□□□□□□□□□□□□□□Auxiliary Table□□□□□□□□□□□□□□□Auxiliary Table□□□□□□FTS Index Cache□□□□word□□□□□□Auxiliary Table□□□□□□□□□□□□□□merge□□□□□□□□□□□□□□Insert Buffer□□□□□□□□□□Insert Buffer□□□□□□□□□□□□□□□B+□□□□□□□□□FTS Index Cache□□□□□□Insert Buffer□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Auxiliary Table□□□□□□

InnoDB□□□□□□□□□□□□□□□□□□□Auxiliary Table□□□□□□□□□□□□□□□□□□□innodb_ft_aux_table□□□□□□□□□□Auxiliary Table□□□□SQL□□□□□□□test□□□□□fts_a□Auxiliary Table□

```
mysql□SET GLOBAL innodb_ft_aux_table='test/fts_a';

Query OK,0 rows affected(0.00 sec)
```

□□□□□□□□□□□□□□□□□□□□information_schema□□□□□□□□INNODB_FT_INDEX_TABLE□□□□fts_a□□□□□□□□□□□

□□□□□□□□□Oracle 11g□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□FTS Index Cache□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□FTS Index Cache□□□□□□□□□□□□□□Auxiliary Table□□□□□□□□□□□□□□□□□□□□□□□□FTS Index Cache□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□FTS Index Cache□□

□□innodb_ft_cache_size□□□□□FTS Index Cache□□□□□□□□□32M□□□□□□□□□□□□□□(word□ilist)□□□□□□□□□□□Auxiliary Table□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

FTS Document ID□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□word□□□□□□□InnoDB□□□□□□□□□FTS_DOC_ID□□□□□□□□□BIGINT UNSIGNED NOT NULL□□□InnoDB□□□□□□□□□□□□□□□□□□□FTS_DOC_ID_INDEX□Unique Index□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□FTS_DOC_ID□□□□□□□Unique Index□□□□□□□FTS_DOC_ID□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□

---

```
mysql□CREATE TABLE fts_a(

-□FTS_DOC_ID INT UNSIGNED AUTO_INCREMENT NOT NULL,

-□body TEXT,

-□PRIMARY KEY(FTS_DOC_ID)

-□);

ERROR 1166(42000):Incorrect column name'FTS_DOC_ID'
```

---

□□□□□□□□□□□□□□□□□□FTS_DOC_ID□□□□□INT□□□□BIG INT□□□□□□□□□□□□□□Incorrect column name'FTS_DOC_ID'□□□□□□□□□□□□□□□□□□□□□□□□

---

```
mysql□CREATE TABLE fts_a(
```

```
-　FTS_DOC_ID BIGINT UNSIGNED AUTO_INCREMENT NOT NULL,

-　body TEXT,

-　PRIMARY KEY(FTS_DOC_ID)

-　);

Query OK,0 rows affected(0.02 sec)
```

--------------------------------------------------

当对全文检索进行查询操作时，索引缓存会批量刷新到Auxiliary
Table中。但是当前将删除操作的FTS Cache Index并不会删除Auxiliary Table中
的记录，因此InnoDB存储引擎会记录FTS Document ID，并将其保存在DELETED
auxiliary table中。可以通过参数innodb_ft_aux_table来观察，例如通过查询表
information_schema架构下的INNODB_FT_DELETED来检测已删除的FTS
Document ID。

随着每次对DML操作时全文索引缓存的变化，会有越来越多的DELETED记录，故会
导致全文检索的查询会变慢。当删除的记录达到一定数量后，用户可以手动地让
InnoDB存储引擎来删除记录，以此提高全文检索的查询性能。用户可以通过命令
OPTIMIZE TABLE。由于OPTIMIZE TABLE还会进行一些其他的操作，如
Cardinality的重新计算，若只希望对倒排索引进行操作，那么可以通过参数
innodb_optimize_fulltext_only进行设置：

--------------------------------------------------

```
mysql　SET GLOBAL innodb_optimize_fulltext_only=1;

mysql　OPTIMIZE TABLEfts_a;
```

--------------------------------------------------

由于对全文检索进行的删除操作OPTIMIZE TABLE可能需要较长的时间，为此可以通
过参数限制每次删除的分词数量，通过参数
innodb_ft_num_word_optimize控制，该参数的默认值设置的值为
2000。

现在来看如何在表中使用全文检索。首先创建表fts_a：

--------------------------------------------------

```
CREATE TABLE fts_a(

FTS_DOC_ID BIGINT UNSIGNED AUTO_INCREMENT NOT NULL,

body TEXT,
```

```
PRIMARY KEY(FTS_DOC_ID)

);

INSERT INTO fts_a

SELECT NULL,'Pease porridge in the pot';

INSERT INTO fts_a

SELECT NULL,'Pease porridge hot,pease porridge cold';

INSERT INTO fts_a

SELECT NULL,'Nine days old';

INSERT INTO fts_a

SELECT NULL,'Some like it hot,some like it cold';

INSERT INTO fts_a

SELECT NULL,'Some like it in the pot';

INSERT INTO fts_a

SELECT NULL,'Nine days old';

INSERT INTO fts_a

SELECT NULL,'I like code days';

CREATE FULLTEXT INDEX idx_fts ON fts_a(body);
```

---

# 上述代码创建表fts_a，其中body列为需要进行全文检索的字段。通过命令创建FULLTEXT索引。由于之前已经在表上显式创建了名为FTS_DOC_ID的列，因此数据库会自动将该列作为全文检索的文档标识。下面来具体看表fts_a中的数据：

---

```
mysql□SELECT*FROM fts_a;

+------------+----------------------------------------+

|FTS_DOC_ID|body|

+------------+----------------------------------------+

|1|Pease porridge in the pot|

|2|Pease porridge hot,pease porridge cold|

|3|Nine days old|

|4|Some like it hot,some like it cold|
```

|5|Some like it in the pot|

|6|Nine days old|

|7|I like code days|

+------------+----------------------------------------+

7 rows in set(0.00 sec)

# 此时我们需要innodb_ft_aux_table来查看分词对应的信息

mysql：SET GLOBAL innodb_ft_aux_table='test/fts_a';

Query OK,0 rows affected(0.00 sec)

mysql：SELECT*FROM information_schema.INNODB_FT_INDEX_TABLE;

+----------+--------------+-------------+-----------+--------+----------+

|WORD|FIRST_DOC_ID|LAST_DOC_ID|DOC_COUNT|DOC_ID|POSITION|

+----------+--------------+-------------+-----------+--------+----------+

|code|7|7|1|7|7|

|cold|2|4|2|2|35|

|cold|2|4|2|4|31|

|days|3|7|3|3|5|

|days|3|7|3|6|5|

|days|3|7|3|7|12|

|hot|2|4|2|2|15|

|hot|2|4|2|4|13|

|like|4|7|3|4|5|

|like|4|7|3|4|18|

|like|4|7|3|5|5|

|like|4|7|3|7|2|

|nine|3|6|2|3|0|

|nine|3|6|2|6|0|

|old|3|6|2|3|10|

|old|3|6|2|6|10|

|pease|1|2|2|1|0|

|pease|1|2|2|2|0|

|pease|1|2|2|2|20|

|porridge|1|2|2|1|6|

|porridge|1|2|2|2|6|

|porridge|1|2|2|2|20|

|pot|1|5|2|1|22|

|pot|1|5|2|5|20|

|some|4|5|2|4|0|

|some|4|5|2|4|18|

|some|4|5|2|5|0|

+----------+--------------+-------------+------------+--------+----------+

27 rows in set(0.00 sec)

---

# 列名中除了word字段外还有DOC_ID、POSITION信息，还有新的FIRST_DOC_ID、LAST_DOC_ID以及DOC_COUNT字段。这些字段表示word第一次出现的文档ID，最后一次出现的文档ID，以及word在多少个文档中存在。

如果用户删除了表中具有FTS_DOC_ID为7的记录：

---

mysql\DELETE FROM test.fts_a WHERE FTS_DOC_ID=7;

Query OK,1 row affected(0.00 sec)

---

# 那么用户可以通过查询InnoDB存储引擎内部的另一张表来得到被删除的文档ID。在下面的例子中，用户可以看到被删除的文档ID为7，表DELETED中记录了删除文档的文档ID。如下所示：

---

mysql\SELECT*FROM INNODB_FT_DELETED;

+--------+

|DOC_ID|

+--------+

|7|

+--------+

1 row in set(0.00 sec)

---

## 这样做的好处就是将ID号最小的一条INNODB_FT_DELETED记录进行删除操作，同时相应的全文索引信息也会被删除。这时再来看一下SQL的结果：

mysql＞SET GLOBAL innodb_optimize_fulltext_only=1;

Query OK,0 rows affected(0.00 sec)

mysql＞OPTIMIZE TABLE test.fts_a;

+------------+----------+----------+----------+

|Table|Op|Msg_type|Msg_text|

+------------+----------+----------+----------+

|test.fts_a|optimize|status|OK|

+------------+----------+----------+----------+

1 row in set(0.01 sec)

mysql＞SELECT*FROM INNODB_FT_DELETED;

+--------+

|DOC_ID|

+--------+

|7|

+--------+

1 row in set(0.00 sec)

mysql＞SELECT*FROM INNODB_FT_BEING_DELETED;

+--------+

|DOC_ID|

+--------+

|7|

+--------+

1 row in set(0.00 sec)

---

在用户用OPTIMIZE TABLE命令删除已经被标记删除的记录，若设置的被删除的ID都会记录到INNODB_FT_BEING_DELETED。上述插入7这条记录ID虽然是递增的，但是由于它指定了ID，故InnoDB存储引擎判断其非法，抛出异常。

---

mysql＞INSERT INTO test.fts_a SELECT 7,'I like this days';

ERROR 182(HY000):Invalid InnoDB FTS Doc ID

---

stopword列表（stopword list）表示该列表中的word不需要对其进行索引分词操作。例如，对于the这个单词，由于其不具有具体的意义，因此将其视为stopword。InnoDB存储引擎有一张默认的stopword列表，其在information_schema架构下，表名为INNODB_FT_DEFAULT_STOPWORD，默认共有36个stopword。此外用户也可以通过参数innodb_ft_server_stopword_table来自定义stopword列表。如：

---

mysql＞CREATE TABLE user_stopword(

-＞value VARCHAR(30)

-＞)ENGINE=INNODB;

Query OK,0 rows affected(0.03 sec)

mysql＞SET GLOBAL

-＞innodb_ft_server_stopword_table="test/user_stopword";

Query OK,0 rows affected(0.00 sec)

---

当前InnoDB存储引擎的全文检索还存在以下的限制：

□每张表只能有一个全文检索的索引。

□由多列组合而成的全文检索的索引列必须使用相同的字符集与排序规则。

□不支持没有单词界定符（delimiter）的语言，如中文、日语、韩语等。

# 5.8.4　全文检索

MySQL数据库支持全文检索（Full-Text Search）的查询，其语法为：

```
MATCH(col1,col2,...)AGAINST(expr[search_modifier])

search_modifier:

{

IN NATURAL LANGUAGE MODE

|IN NATURAL LANGUAGE MODE WITH QUERY EXPANSION

|IN BOOLEAN MODE

|WITH QUERY EXPANSION

}
```

MySQL数据库通过MATCH()...AGAINST()函数来进行全文检索的查询。MATCH指定了需要被查询的列，AGAINST指定了使用何种方法去进行查询。下面将对各种查询模式进行介绍。

## 1.Natural Language

全文检索通过MATCH函数进行查询，默认采用Natural Language模式，其表示查询带有指定word的文档。对于5.8.3小节中创建的表fts_a，查询body字段中带有Pease的文档，若不使用全文检索技术，则允许使用下述SQL语句：

```
mysql〉SELECT*FROM fts_a WHERE body LIKE'%Pease%';
```

由于不能使用B+树索引，因此上述SQL语句需要对整个表进行扫描以找到结果。而利用全文检索技术，可以用下述SQL语句进行查询：

```
mysql〉SELECT*FROM fts_a

-〉WHERE MATCH(body)

-〉AGAINST('Porridge'IN NATURAL LANGUAGE MODE);

+------------+----------------------------------------+
```

|FTS_DOC_ID|body|

+------------+----------------------------------------+

|2|Pease porridge hot,pease porridge cold|

|1|Pease porridge in the pot|

+------------+----------------------------------------+

2 rows in set(0.00 sec)

# 通过NATURAL LANGUAGE MODE的全文检索是默认的全文检索模式。下面将显示通过自然语言查询SQL语句的全文检索。

SELECT*FROM fts_a WHERE MATCH(body)AGAINST('Porridge');

# 再通过查看SQL优化器来观察全文检索。

mysql〉EXPLAIN SELECT*FROM fts_a

-〉WHERE MATCH(body)AGAINST('Porridge')\G;

***************************1.row**************************

id:1

select_type:SIMPLE

table:fts_a

type:fulltext

possible_keys:idx_fts

key:idx_fts

key_len:0

ref:NULL

rows:1

Extra:Using where

1 row in set(0.00 sec)

我们可以发现type这里为fulltext，即使用全文检索的倒排索引，key这里显示的
idx_fts。倒排索引会把所有分出来的单词都放到一个索引里，用户进行查询时，根据
MATCH直接就可以进行查询。

---

mysql＞SELECT*FROM fts_b

-＞WHERE MATCH(body)AGAINST('Porridge');

ERROR 1191(HY000):Can't find FULLTEXT index matching the column list

---

在WHERE条件中使用MATCH函数，查询返回的结果是根据相关性（Relevance）进行降
序排序的，即相关性最高的结果放在第一位。相关性的值是一个非负浮点数字，0表示没有
任何的相关性。MySQL数据库使用如下公式进行相关性的统计计算：

□word是否在文档中出现

□word在文档中出现的次数

□word在索引列中的数量

□多少个文档包含该word

上述查询的结果中Porridge分别2次出现在数据库表的两行记录中，因此在全文检索下，
可以发现使用MATCH函数并没有将其过滤，下面的SQL语句：

---

mysql＞SELECT count(*)

-＞FROM fts_a WHERE

-＞MATCH(body)AGAINST('Porridge'IN NATURAL LANGUAGE MODE);

+-------------------+

|count(FTS_DOC_ID)|

+-------------------+

|2|

+-------------------+

1 row in set(0.00 sec)

## 上述SQL语句的计算结果如下。

```
mysql>SELECT

->COUNT(IF(MATCH(body)

->AGAINST('Porridge'IN NATURAL LANGUAGE MODE),1,NULL))

->AS count

->FROM fts_a;

+-------+

|count|

+-------+

|2|

+-------+

1 row in set(0.00 sec)
```

## 上面这条SQL语句返回文档中匹配的记录数量。可以看到通过上述SQL语句的确可以获得相应的结果，然而这条SQL语句相当得不友好。其实用户更希望获得的结果是具体每个文档的相关性。这可以将相关性进行排序。下面的SQL语句用于计算

## 文档的相关性，通过这条SQL语句可以得到期望的结果。

```
mysql>SELECT fts_doc_id,body,

->MATCH(body)AGAINST('Porridge'IN NATURAL LANGUAGE MODE)

->AS Relevance

->FROM fts_a;

+------------+-----------------------------------------+-------------------+

|fts_doc_id|body|Relevance|

+------------+-----------------------------------------+-------------------+

|1|Pease porridge in the pot|0.2960100471973419|

|2|Pease porridge hot,pease porridge cold|0.5920200943946838|

|3|Nine days old|0|
```

|4|Some like it hot,some like it cold|0|

|5|Some like it in the pot|0|

|6|Nine days old|0|

|7|I like hot and code days|0|

+------------+----------------------------------------+--------------------+

7 rows in set(0.01 sec)

---

通过InnoDB存储引擎的全文检索可以查询各表中的

单词（word）、stopword、单词所在的文档数量等信息

其中单词（word）必须是在区间[innodb_ft_min_token_size和
innodb_ft_max_token_size]之间

我们对于stopword，可以观察一个简单的例子对于the这个单词，我们可以看到它是一个

---

```
mysql　SELECT fts_doc_id AS id,body,

-　MATCH(body)AGAINST('the'IN NATURAL LANGUAGE MODE)

-　AS rl

-　FROM fts_a;
```

+----+----------------------------------------+------+

|id|body|rl|

+----+----------------------------------------+------+

|1|Pease porridge in the pot|0|

|2|Pease porridge hot,pease porridge cold|0|

|3|Nine days old|0|

|4|Some like it hot,some like it cold|0|

|5|Some like it in the pot|0|

|6|Nine days old|0|

|7|I like hot and code days|0|

+----+----------------------------------------+------+

---

可以发现，单词the在文档1和5中出现了，但是stopword，因此其出现次数为0。

参数innodb_ft_min_token_size和innodb_ft_max_token_size控制InnoDB存储引擎查询字符的长度，当长度小于innodb_ft_min_token_size，或者长度大于innodb_ft_max_token_size时，会忽略该词的搜索。在InnoDB存储引擎中，innodb_ft_min_token_size的默认值为3，innodb_ft_max_token_size的默认值为84。

## 2.Boolean

MySQL数据库支持全文检索IN BOOLEAN MODE的查询，这称为布尔全文检索。它和之前介绍的全文检索区别在于，要查询的字符串带有一个或多个操作符，如查询Pease但不包括hot，可以设置为+、-等，下面简单介绍各个操作符的含义。

---

```
mysql＞SELECT*FROM fts_a

-＞WHERE MATCH(body)AGAINST('+Pease-hot'IN BOOLEAN MODE)\G;

***************************1.row***************************

FTS_DOC_ID:1

body:Pease porridge in the pot
```

---

Boolean全文检索支持以下几种操作符：

□+表示该word必须存在。

□-表示该word必须被排除。

□（no operator）表示该word是可选的，但是如果出现，其相关性会更高。

□@distance表示查询的多个单词之间的距离是否在distance之内，distance的单位是字节。这种全文检索的查询也称为Proximity Search。如MATCH（body）AGAINST（'"Pease pot"@30'IN BOOLEAN MODE）表示字符Pease和pot之间的距离需在30字节内。

□口口口口口口口口口口口口口口口

□口口口口口口口口口口口口口口口

□口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口

□*口口口口口口口口口口口口lik*口口口口口口口lik口like口口口口likes口

□"口口口口口

口口口口口口口口口口口口口口口口fts_a口口口口口口口口口口口口SQL口口口口口pease口口hot口口口口

---

```
mysql口SELECT*FROM fts_a

-口WHERE MATCH(body)AGAINST('+Pease+hot'IN BOOLEAN MODE)\G;

***************************1.row***************************

FTS_DOC_ID:2

body:Pease porridge hot,pease porridge cold

1 row in set(0.00 sec)
```

---

## 口口口SQL口口口口口pease口口口hot口口口口

---

```
mysql口SELECT*FROM fts_a

-口WHERE MATCH(body)AGAINST('+Pease-hot'IN BOOLEAN MODE)\G;

***************************1.row***************************

FTS_DOC_ID:1

body:Pease porridge in the pot

1 row in set(0.00 sec)
```

---

## 口口口SQL口口口口口pease口口hot口口口口

---

```
mysql口SELECT*FROM fts_a
```

```
-> WHERE MATCH(body)AGAINST('Pease hot'IN BOOLEAN MODE);

+-----------+----------------------------------------+

|FTS_DOC_ID|body|

+-----------+----------------------------------------+

|2|Pease porridge hot,pease porridge cold|

|1|Pease porridge in the pot|

|4|Some like it hot,some like it cold|

|7|I like hot and code days|

+-----------+----------------------------------------+

4 rows in set(0.00 sec)
```

## 邻近度SQL查询语句（Proximity Search）

```
mysql>SELECT fts_doc_id,body FROM fts_a

-> WHERE MATCH(body)

-> AGAINST('"Pease pot"@30'IN BOOLEAN MODE)\G;

***************************1.row***************************

fts_doc_id:1

body:Pease porridge in the pot

1 row in set(0.01 sec)

mysql>SELECT fts_doc_id,body FROM fts_a

-> WHERE MATCH(body)

-> AGAINST('"Pease pot"@10'IN BOOLEAN MODE);

Empty set(0.01 sec)
```

## 上述查询文档1中由于Pease和pot单词的间隔为22，因此第一个查询@30可以找到文档，而第二个查询@10则不能找到文档。

```
mysql>SELECT fts_doc_id,body,

-> MATCH(body)AGAINST('like pot'IN BOOLEAN MODE)
```

-　AS Relevance FROM fts_a;

+------------+------------------------------------------+--------------------+

|fts_doc_id|body|Relevance|

+------------+------------------------------------------+--------------------+

|1|Pease porridge in the pot|1.2960100173950195|

|2|Pease porridge hot,pease porridge cold|0|

|3|Nine days old|0|

|4|Some like it hot,some like it cold|0.27081382274627686|

|5|Some like it in the pot|1.4314169883728027|

|6|Nine days old|0|

|7|I like hot and code days|0.13540691137313843|

+------------+------------------------------------------+--------------------+

7 rows in set(0.00 sec)

---

# 上述SQL语句查询根据字段中包含like和pot的文档,虽然结果没有字段中包含pot这个单词的文档,而第4个文档包含like但是不包含pot,所以也只有1和5这两

# 个文档。再来看包含"-some"这个单词的文档,具体如下所示:

---

mysql　SELECT fts_doc_id,body,

-　MATCH(body)AGAINST('like　hot　some'IN BOOLEAN MODE)

-　AS Relevance

-　FROM fts_a;

+------------+------------------------------------------+--------------------+

|fts_doc_id|body|Relevance|

+------------+------------------------------------------+--------------------+

|1|Pease porridge in the pot|0|

|2|Pease porridge hot,pease porridge cold|1.2960100173950195|

|3|Nine days old|0|

|4|Some like it hot,some like it cold|1.158843994140625|

|5|Some like it in the pot|-0.5685830116271973|

|6|Nine days old|0|

|7|I like hot and code days|0.13540691137313843|

+------------+----------------------------------------+--------------------+

7 rows in set(0.00 sec)

上述查询结果中5这个文档展示了评分负的情况，由于like是该文档中存在，而some却不是其中一个单词，所以计算得分为负。

接着看如下SQL查询：

mysql〉SELECT*FROM fts_a

-〉WHERE MATCH(body)AGAINST('po*'IN BOOLEAN MODE);

+------------+----------------------------------------+

|FTS_DOC_ID|body|

+------------+----------------------------------------+

|2|Pease porridge hot,pease porridge cold|

|1|Pease porridge in the pot|

|5|Some like it in the pot|

+------------+----------------------------------------+

3 rows in set(0.00 sec)

上述查询找出了所有包含以po开头的单词，如porridge和pot。

再来看下面的SQL查询结果：

mysql〉SELECT*FROM fts_a

-〉WHERE MATCH(body)AGAINST('like hot'IN BOOLEAN MODE);

+------------+----------------------------------------+

|FTS_DOC_ID|body|

+------------+----------------------------------------+

```
|4|Some like it hot,some like it cold|

|7|I like hot and code days|

|2|Pease porridge hot,pease porridge cold|

|5|Some like it in the pot|

+------------+-------------------------------------+

4 rows in set(0.00 sec)

mysql□SELECT*FROM fts_a

-□WHERE MATCH(body)AGAINST('"like hot"'IN BOOLEAN MODE);

+------------+-------------------------+

|FTS_DOC_ID|body|

+------------+-------------------------+

|7|I like hot and code days|

+------------+-------------------------+

1 row in set(0.00 sec)
```

上述例子中第一条SQL语句由于没有""号，即表示查询的结果中存在字符串或者的文档，故有4个结果。而第二条SQL由于有"like hot"，故表示必须相连的like hot字符串，因此只有4条记录返回。

## 3.Query Expansion

MySQL数据库支持全文检索的扩展查询。这种查询通常在查询的关键词太短，用户需要□implied knowledge□（隐含知识）时进行。例如，对于单词database的查询，用户可能希望查询的不仅仅是包含database的文档，可能还指那些包含MySQL、Oracle、DB2、RDBMS的单词。而这时可以使用Query Expansion模式来开启全文检索的□implied knowledge□。

通过在查询短语中添加WITH QUERY EXPANSION□IN NATURAL LANGUAGE MODE WITH QUERY EXPANSION□可以开启□blind query expansion□（又称□automatic relevance feedback□）。该查询分为两个阶段。

□第一阶段：根据搜索的单词进行全文索引查询。

□第二阶段：根据第一阶段产生的分词再进行一次全文检索的查询。

接着来看一个具体的例子，还是对于表articles，

```
CREATE TABLE articles(

id INT UNSIGNED AUTO_INCREMENT NOT NULL PRIMARY KEY,

title VARCHAR(200),

body TEXT,

FULLTEXT(title,body)

)ENGINE=InnoDB;

INSERT INTO articles(title,body)VALUES

('MySQL Tutorial','DBMS stands for DataBase...'),

('How To Use MySQL Well','After you went through a...'),

('Optimizing MySQL','In this tutorial we will show...'),

('1001 MySQL Tricks','1.Never run mysqld as root.2....'),

('MySQL vs.YourSQL','In the following database comparison...'),

('MySQL Security','When configured properly,MySQL...'),

('Tuning DB2','For IBM database...'),

('IBM History','DB2 hitory for IBM...');
```

## 上面这段代码创建了一个包含FTS_DOC_ID列的InnoDB全文索引表,下面通过一个实例来说明。在articles表中查找包含关键词title、body以及包含关键词database的文章,示例代码如下:

```
mysql□SELECT*FROM articles

-□WHERE MATCH(title,body)

-□AGAINST('database'IN NATURAL LANGUAGE MODE);

+----+-----------------+----------------------------------------+

|id|title|body|

+----+-----------------+----------------------------------------+

|1|MySQL Tutorial|DBMS stands for DataBase...|

|5|MySQL vs.YourSQL|In the following database comparison...|

|7|Tuning DB2|For IBM database...|
```

```
+----+------------------+-----------------------------------------+
```

3 rows in set(0.00 sec)

上面的例子返回了3个结果（body字段包含database的结果）。如果采用Query Expansion，搜索结果将会变成如下所示：

```
mysql>SELECT*FROM articles

->WHERE MATCH(title,body)

->AGAINST('database'WITH QUERY EXPANSION);
```

```
+----+---------------------+-----------------------------------------+

|id|title|body|

+----+---------------------+-----------------------------------------+

|5|MySQL vs.YourSQL|In the following database comparison...|

|1|MySQL Tutorial|DBMS stands for DataBase...|

|7|Tuning DB2|For IBM database...|

|8|IBM History|DB2 hitory for IBM...|

|3|Optimizing MySQL|In this tutorial we will show...|

|6|MySQL Security|When configured properly,MySQL...|

|2|How To Use MySQL Well|After you went through a...|

|4|1001 MySQL Tricks|1.Never run mysqld as root.2....|

+----+---------------------+-----------------------------------------+
```

8 rows in set(0.00 sec)

可以看到，返回了8条结果，其中包含单词database的行排在前列，同时也包含title或body中包含MySQL、DB2的文章。这就是Query Expansion。

采用Query Expansion可以使搜索结果更加精确，这是由于它在原来搜索的单词基础上，增加了新的可能相关的词汇。

# 5.9 小结

本章介绍了表的物理存储、表的逻辑存储以及索引的内部实现和B+树索引的使用。与其他数据库一样，InnoDB1.2版本支持的索引类型也非常广泛，这有助于数据库性能的提高。本章详细介绍了B+树索引，以及如何通过命令来查看索引的使用和内部的存储状态。

# 第6章 锁

开发多用户、数据库驱动的应用时，最大的一个难点是：一方面要最大程度地利用数据库的并发访问，另一方面还要确保每个用户能以一致的方式读取和修改数据。为此就有了锁（locking）的机制，同时这也是数据库系统区别于文件系统的一个关键特性。对于InnoDB存储引擎而言，MySQL数据库其本身对于锁的支持是不尽相同的，一些存储引擎并不支持，如Oracle中的共享锁。因此本章将主要讨论存储引擎——InnoDB的锁实现的一些内部机制。

很多人对各种数据库的InnoDB存储引擎中的锁有误解，或者直接将InnoDB存储引擎的锁的实现和其他数据库，如MyISAM、Oracle、SQL Server相混淆。这种理解偏差会导致用户错误地使用"锁"，这反而不利于应用程序的正常运行，影响应用程序并发性的同时，也会使你对InnoDB存储引擎产生怀疑，觉得其根本不支持锁或者锁的功能存在问题。

# 6.1 □□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ [1]□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□LRU□□□□□□□□□□□□□LRU□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□MySQL□□□□MyISAM□NDB Cluster□□□□□□□□□□MySQL□□□□□□□□□□□□□□Microsoft SQL Server□Oracle□□□□□□□□□□□□□□□□□□□□□□□

□□MyISAM□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□"□□"□MyISAM□□□□□□□□□□□□□□□□□□□□□□□□□Microsoft SQL Server□□□□□Microsoft SQL Server 2005□□□□□□□□□□□□□□□□□□□□MyISAM□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□2005□□□Microsoft SQL Server□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□Microsoft SQL Server□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

InnoDB□□□□□□□□□□□□Oracle□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

[1] □□□□□□□□□□"□□□□"□□□□□□□"□□□"□

## 6.2 lock□latch

□□□□□□□□□□□□□□□□□□lock□latch□□□□□□□□□□lock□latch□□□□□□□□"□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□lock□

latch□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□latch□□□□□□□mutex□□□□□□□□□rwlock□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

lock□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□lock□□□□□□□□□□commit□rollback□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□lock□□□□□□□□□□□□□□□□□□□□□□□□□□6-1□□□lock□latch□□□□□

表 6-1　lock 与 latch 的比较

|  | lock | latch |
|---|---|---|
| 对象 | 事务 | 线程 |
| 保护 | 数据库内容 | 内存数据结构 |
| 持续时间 | 整个事务过程 | 临界资源 |
| 模式 | 行锁、表锁、意向锁 | 读写锁、互斥量 |
| 死锁 | 通过 waits-for graph、time out 等机制进行死锁检测与处理 | 无死锁检测与处理机制。仅通过应用程序加锁的顺序（lock leveling）保证无死锁的情况发生 |
| 存在于 | Lock Manager 的哈希表中 | 每个数据结构的对象中 |

查看InnoDB存储引擎中的latch可以通过命令SHOW ENGINE INNODB MUTEX，如图6-1所示。

```
mysql> SHOW ENGINE INNODB MUTEX;
+--------+----------------+-------------+
| Type   | Name           | Status      |
+--------+----------------+-------------+
| InnoDB | srv0srv.c:1020 | os_waits=5  |
| InnoDB | log0log.c:833  | os_waits=3  |
+--------+----------------+-------------+
2 rows in set (0.03 sec)
```

图 6-1　非Debug版本下SHOW ENGINE INNODB MUTEX的latch

在Debug版本下，通过命令SHOW ENGINE INNODB MUTEX可以看到latch的更多信息，如图6-2所示。

```
mysql> SHOW ENGINE INNODB MUTEX;
+--------+-----------------------+-------------------------------------------------------------------------------+
| Type   | Name                  | Status                                                                        |
+--------+-----------------------+-------------------------------------------------------------------------------+
| InnoDB | &kernel_mutex:srv0srv.c | count=54, spin_waits=6, spin_rounds=60, os_waits=3, os_yields=3, os_wait_times=0 |
| InnoDB | log0log.c:833         | os_waits=2                                                                    |
| InnoDB | rw_lock_mutexes       | count=0, spin_waits=0, spin_rounds=0, os_waits=0, os_yields=0, os_wait_times=0 |
+--------+-----------------------+-------------------------------------------------------------------------------+
3 rows in set (0.01 sec)
```

表 6-2 在Debug版本下可以看到latch

的更多信息，如表所示。其中Type表示是否为InnoDB，Name表示打开latch时源代码的文件名，以及所在行数，Status表示在非Debug版本下，显示的os_waits，即表中的count、spin_waits、spin_rounds、os_yields、os_wait_times，这些变量的说明见表6-2。

表6-2 命令 SHOW ENGINE INNODB MUTEX 输出结果说明

| 名称 | 说明 |
|---|---|
| count | mutex 被请求的次数 |
| spin_waits | spin lock（自旋锁）的次数，InnoDB 存储引擎 latch 在不能获得锁时首先进行自旋，若自旋后还不能获得锁，则进入等待状态 |
| spin_rounds | 自旋内部循环的总次数，每次自旋的内部循环是一个随机数。spin_rounds/spain_waits 表示平均每次自旋所需的内部循环次数 |
| os_waits | 表示操作系统等待的次数。当 spin lock 通过自旋还不能获得 latch 时，则会进入操作系统等待状态，等待被唤醒 |
| os_yields | 进行 os_thread_yield 唤醒操作的次数 |
| os_wait_times | 操作系统等待的时间，单位是 ms |

相对于锁来说，读者需要了解的仅是如何查看当前数据库中锁的请求，然后判断由此

关于latch以及相关lock的信息，还可以通过如下命令进行查看：SHOW ENGINE INNODB STATUS、information_schema架构下的INNODB_TRX、INNODB_LOCKS、INNODB_LOCK_WAITS。相关内容将在下一小节中进行介绍。

# 6.3　InnoDB□□□□□□□□

## 6.3.1　□□□□

InnoDB□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□S Lock□□□□□□□□□□□□□□□□

□□□□□□X Lock□□□□□□□□□□□□□□□□□□□

□□□□□□□□T1□□□□□□□r□□□□□□□□□□□□□T2□□□□□□□□r□□□□□□□□□□□□□□□□□□□□□r□□□□□□□□□□□□□□□Lock Compatible□□□□□□□□□□□□□□□□T3□□□□r□□□□□□□□□□□□□□□□□T1□T2□□□□r□□□□□□——□□□□□□□□□□□□□□□□6-3□□□□□□□□□□□□□□□□□□□

表6-3　排他锁和共享锁的兼容性

|  | X | S |
|---|---|---|
| X | 不兼容 | 不兼容 |
| S | 不兼容 | 兼容 |

□□□6-3□□□□□□X□□□□□□□□□□□□□□□□□□S□□□□S□□□□□□□□□□□□□□□□□□□□□□S□X□□□□□□□□□□□□□□□□□□□□□□row□□□□□□□□□□□□□

□□□□InnoDB□□□□□□□□□□□□□granular□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Intention Lock□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□fine granularity□□□□□□□□□□□□6-3□□□□□

## 图 6-3　意向锁

为了允许行锁和表锁共存，实现多粒度锁机制，InnoDB 还有两种内部使用的意向锁，这两种意向锁都是表锁，如图 6-3 所示。当事务要获取某个 r 行的 X 锁（排他锁）时（图 A），首先要获取表的 IX 锁（意向排他锁），接着再获取 r 的 X 锁；当事务要获取某个 1 行的 S 锁时（获取 1 行的共享 S 锁），首先要获取表的 S 锁，接着再获取 r 行的 1 行的 IX 锁。至于为什么要这样设计，读者可以简单地认为这样设计较为方便。

InnoDB 存储引擎支持多粒度锁定，这种锁定允许事务在行级上的锁和表级上的锁同时存在。为了支持在不同粒度上进行加锁操作，

1）意向共享锁（IS Lock），事务想要获得一张表中某几行的共享锁。

2）意向排他锁（IX Lock），事务想要获得一张表中某几行的排他锁。

由于 InnoDB 存储引擎支持的是行级别的锁，因此意向锁其实不会阻塞除全表扫以外的任何请求。故表级意向锁与行级锁的兼容性如表 6-4 所示。

## 表 6-4　InnoDB 存储引擎中锁的兼容性

|  | IS | IX | S | X |
|---|---|---|---|---|
| IS | 兼容 | 兼容 | 兼容 | 不兼容 |
| IX | 兼容 | 兼容 | 不兼容 | 不兼容 |
| S | 兼容 | 不兼容 | 兼容 | 不兼容 |
| X | 不兼容 | 不兼容 | 不兼容 | 不兼容 |

用户可以通过以下 SHOW ENGINE INNODB STATUS 命令来查看当前锁请求的信息：

mysql　SHOW ENGINE INNODB STATUS\G;

……

------------

TRANSACTIONS

------------

Trx id counter 48B89BF

Purge done for trx's n:o　48B89BA undo n:o　0

History list length 0

LIST OF TRANSACTIONS FOR EACH SESSION:

---TRANSACTION 0,not started,process no 13757,OS thread id 1255176512

MySQL thread id 42,query id 80424887 localhost root

show engine innodb status

---TRANSACTION 48B89BE,ACTIVE 193 sec,process no 13757,OS thread id 1254910272 starting index read

mysql tables in use 1,locked 1

LOCK WAIT 2 lock struct(s),heap size 368,1 row lock(s)

MySQL thread id 41,query id 80424886 localhost root Sending data

select*from t where a　4 lock in share mode

-------TRX HAS BEEN WAITING 2 SEC FOR THIS LOCK TO BE GRANTED:

RECORD LOCKS space id 30 page no 3 n bits 72 index'PRIMARY'of table'test'.'t'trx id 48B89BE lock mode S waiting

------------------

TABLE LOCK table'test'.'t'trx id 48B89BE lock mode IS

RECORD LOCKS space id 30 page no 3 n bits 72 index'PRIMARY'of table'test'.'t'trx id 48B89BE lock mode S waiting

---TRANSACTION 48B89BD,ACTIVE 205 sec,process no 13757,OS thread id 1257838912

2 lock struct(s),heap size 368,1 row lock(s)

MySQL thread id 40,query id 80424881 localhost root

TABLE LOCK table'test'.'t'trx id 48B89BD lock mode IX

RECORD LOCKS space id 30 page no 3 n bits 72 index'PRIMARY'of table'test'.'t'trx id 48B89BD lock_mode X locks rec but not gap

----------------------------

END OF INNODB MONITOR OUTPUT

```
==========================

1 row in set(0.01 sec)
```

─────────────────────────

对应的SQL语句select*from t where a＝4 lock in share mode，产生了RECORD LOCKS space id 30 page no 3 n bits 72 index'PRIMARY'of table'test'.'t'trx id 48B89BD lock_mode X locks rec but not gap的加锁，其中同样是locks rec but not gap，即在行锁上加上了一个排他锁。

在InnoDB 1.0版本之前，用户只能通过命令SHOW FULL PROCESSLIST、SHOW ENGINE INNODB STATUS等来查看当前数据库中锁的请求，然后再判断事务锁的情况。从InnoDB1.0开始，在INFORMATION_SCHEMA架构下添加了INNODB_TRX、INNODB_LOCKS、INNODB_LOCK_WAITS。通过这三张表，用户可以更简单地监控当前事务并分析可能存在的锁问题。下面将具体分析这三张表。表6-5显示INNODB_TRX表的结构，共有8个字段。

## 表 6-5　表 INNODB_TRX 的结构说明

| 字段名 | 说明 |
| --- | --- |
| trx_id | InnoDB 存储引擎内部唯一的事务 ID |
| trx_state | 当前事务的状态 |
| trx_started | 事务的开始时间 |
| trx_requested_lock_id | 等待事务的锁 ID。如 trx_state 的状态为 LOCK WAIT，那么该值代表当前的事务等待之前事务占用锁资源的 ID。若 trx_state 不是 LOCK WAIT，则该值为 NULL |
| trx_wait_started | 事务等待开始的时间 |
| trx_weight | 事务的权重，反映了一个事务修改和锁住的行数。在 InnoDB 存储引擎中，当发生死锁需要回滚时，InnoDB 存储引擎会选择该值最小的进行回滚 |
| trx_mysql_thread_id | MySQL 中的线程 ID，SHOW PROCESSLIST 显示的结果 |
| trx_query | 事务运行的 SQL 语句 |

## 我们可以通过查看这张表

```
mysql　SELECT*FROM information_schema.INNODB_TRX\G;
***************************1.row***************************
trx_id:7311F4
```

trx_state:LOCK WAIT

trx_started:2010-01-04 10:49:33

trx_requested_lock_id:7311F4:96:3:2

trx_wait_started:2010-01-04 10:49:33

trx_weight:2

trx_mysql_thread_id:471719

trx_query:select*from parent lock in share mode

***************************2.row**************************

trx_id:730FEE

trx_state:RUNNING

trx_started:2010-01-04 10:18:37

trx_requested_lock_id:NULL

trx_wait_started:NULL

trx_weight:2

trx_mysql_thread_id:471718

trx_query:NULL

2 rows in set(0.00 sec)

---

可以从state中发现，在trx_id为730FEE的事务处于运行，而在trx_id为7311F4的事务处于"LOCK WAIT"的状态，而运行的SQL语句为select*from parent lock in share mode。但是这只能告诉我们InnoDB存储引擎当前的状态信息，并不能明确的告诉我们InNODB_LOCKS中具体的锁信息（表6-6）。

## 表6-6 表INNODB_LOCKS的结构

| 字段名 | 说明 |
|---|---|
| lock_id | 锁的ID |
| lock_trx_id | 事务ID |
| lock_mode | 锁的模式 |
| lock_type | 锁的类型，表锁还是行锁 |
| lock_table | 要加锁的表 |
| lock_index | 锁住的索引 |
| lock_space | 锁对象的space id |
| lock_page | 事务锁定页的数量。若是表锁，则该值为NULL |
| lock_rec | 事务锁定行的数量，若是表锁，则该值为NULL |
| lock_data | 事务锁定记录的主键值，若是表锁，则该值为NULL |

## □□□□□□□□□□□□□□INNODB_LOCKS□

```
mysql□SELECT*FROM information_schema.INNODB_LOCKS\G;

***************************1.row***************************

lock_id:7311F4:96:3:2

lock_trx_id:7311F4
```

lock_mode:S

lock_type:RECORD

lock_table:'mytest'.'parent'

lock_index:'PRIMARY'

lock_space:96

lock_page:3

lock_rec:2

lock_data:1

***************************2.row***************************

lock_id:730FEE:96:3:2

lock_trx_id:730FEE

lock_mode:X

lock_type:RECORD

lock_table:'mytest'.'parent'

lock_index:'PRIMARY'

lock_space:96

lock_page:3

lock_rec:2

lock_data:1

2 rows in set(0.00 sec)

从这两条记录我们可以得到，事务trx_id为730FEE的事务正在等待parent表上的X锁，而事务ID为7311F4的事务持有parent表上行锁的S锁，且lock_data都为1。这时如果读者再去查看INNODB_TRX表，就会发现原来的trx_state从"RUNNING"变为了"LOCK WAIT"状态。

通过上述的例子我们可以看到，lock_data这个值并非是"可信"的值。例如当用户运行一个范围查找时，lock_data可能只返回第一行的主键值。与此同时，如果当前资源被锁住了，若锁住的页因为InnoDB存储引擎缓冲池的容量，导致该页从缓冲池中被刷出，则在INNODB_LOCKS表中显示的内容就会为NULL，即InnoDB存储引擎不会从磁盘进行再次查找。

由于在表INNODB_LOCKS中只保存请求当前的事务的等待信息，如果需要查看锁之间的等待关系，还需要对这个表进行自我连接，这样操作比较麻烦。为此，在表INNODB_LOCK_WAITS中记录了锁的等待信息，表INNODB_LOCK_WAITS有4个字段，如表6-7所示。

表6-7  表INNODB_LOCK_WAITS 的结构

| 字段 | 说明 | 字段 | 说明 |
| --- | --- | --- | --- |
| requesting_trx_id | 申请锁资源的事务ID | blocking_trx_id | 阻塞的事务ID |
| requesting_lock_id | 申请的锁的ID | blocking_trx_id | 阻塞的锁的ID |

通过这张表，可以很直观地反映当前事务的等待。如：

```
mysql：SELECT*FROM information_schema.INNODB_LOCK_WAITS\G;

***************************1.row***************************

requesting_trx_id:7311F4

requested_lock_id:7311F4:96:3:2

blocking_trx_id:730FEE

blocking_lock_id:730FEE:96:3:2

1 row in set(0.00 sec)
```

可以通过如下SQL语句将这三张表联合起来，这样就可以比较直观地反映当前运行的事务中哪些事务会产生等待的情况，以及事务ID对应的用户线程，数据库通过将表INNODB_TRX、INNODB_LOCKS、INNODB_LOCK_WAITS联合起来，就可以了解到当前事务在等待什么。如：

```
mysql：SELECT

r.trx_id waiting_trx_id,
```

```
        r.trx_mysql_thread_id waiting_thread,

        r.trx_query waiting_query,

        b.trx_id blocking_trx_id,

        b.trx_mysql_thread_id blocking_thread,

        b.trx_query blocking_query

FROM information_schema.innodb_lock_waits w

INNER JOIN information_schema.innodb_trx b

ON b.trx_id=w.blocking_trx_id

INNER JOIN information_schema.innodb_trx r

ON r.trx_id=w.requesting_trx_id\G;
***************************1.row***************************
waiting_trx_id:73122F

waiting_thread:471719

waiting_query:NULL

blocking_trx_id:7311FC

blocking_thread:471718

blocking_query:NULL

1 row in set(0.00 sec)
```
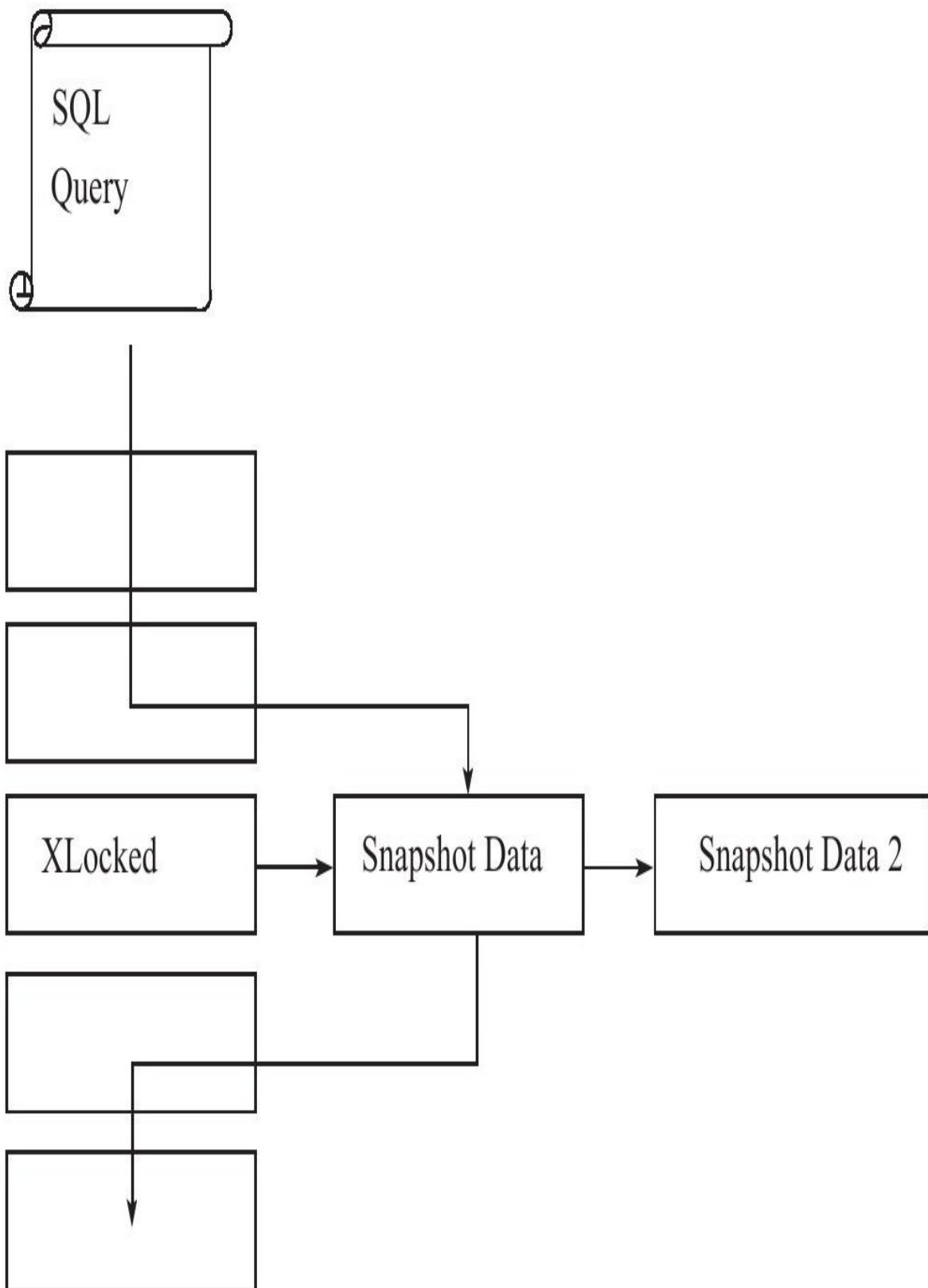
# 6.3.2 　一致性非锁定读

一致性的非锁定读（consistent nonlocking read）是指InnoDB存储引擎通过行多版本控制（multi versioning）的方式来读取当前执行时间数据库中行的数据。如果读取的行正在执行DELETE、UPDATE操作，这时读取操作不会因此去等待行上锁的释放。相反地，InnoDB存储引擎会去读取行的一个快照数据。图6-4直观地

SQL

Query

XLocked → Snapshot Data → Snapshot Data 2

# ❑ 6-4 InnoDB存储引擎中的锁和隔离级别

在6-4节中，我们讨论了InnoDB存储引擎中的锁。我们知道，在执行更新操作时，会对记录加上□X□锁。而查询操作默认是不加锁的，而是通过undo来实现的，undo是如何实现□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□。

□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB存储引擎□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

在本节6-4□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□6-4□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Multi Version Concurrency Control□MVCC□□□

□□□□□□□□□READ COMMITTED□REPEATABLE READ□InnoDB□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□READ COMMITTED□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□REPEATABLE READ□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□A□□□□□SQL□□□□

---

```
#Session A

mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□SELECT*FROM parent WHERE id=1;

+----+

|id|

+----+

|1|

+----+

1 row in set(0.00 sec)
```

---

在□A□□□□□□□□□□□□BEGIN□□□□□□□□□□□□□□□□□parent□id□1□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□B□□□□□□□□□□□□□□□□□□□□□□□B□□□□□□□□

mysql：BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql：UPDATE parent SET id=3 WHERE id=1;

Query OK,1 row affected(0.00 sec)

Rows matched:1 Changed:1 Warnings:0

事务（B）把表parent中id为1的记录修改成id＝3，此时这条被修改的id＝1的记录上会添加X锁。接着事务（A）查询id为1的记录，InnoDB存储引擎事务隔离级别为READ COMMITTED、REPEATETABLE READ时都是采用非锁定的一致性读，故事务（A）能立即得到查询结果。执行的SQL为（SELECT*FROM parent WHERE id=1），观察一下在READ COMMITTED与REPEATABLE READ隔离级别下得到的结果是否相同？

mysql：SELECT*FROM parent WHERE id=1;

+----+

|id|

+----+

|1|

+----+

1 row in set(0.00 sec)

此时得到的id＝1的记录，与最开始的1条记录是一致的，即没有读到已经修改的事务（B）未提交的数据。

#Session B

mysql：commit;

Query OK,0 rows affected(0.01 sec)

此时事务（B）提交完成后，事务（A）再执行SELECT*FROM parent WHERE id=1的SQL语句，在READ COMMITTED与REPEATABLE隔离级别下得到的结果就不同了。在READ COMMITTED隔离级别下，事务会读到已提交事务修改的数据，故最后得到的结果如

可以看到，它使用了fresh snapshot（一个最新的快照），所以可以看到B事务插入的数据。在READ COMMITTED级别下，不存在可重复读的问题。

---

mysql＞SELECT@@tx_isolation\G;

***************************1.row***************************

@@tx_isolation:READ-COMMITTED

1 row in set(0.00 sec)

mysql＞SELECT*FROM parent WHERE id=1;

Empty set(0.00 sec)

---

而对于REPEATABLE READ，因为事务开始时读取的数据一直会维持到结束。故在REPEATABLE READ级别下，不存在不可重复读的问题。

---

mysql＞SELECT@@tx_isolation\G;

***************************1.row***************************

@@tx_isolation:REPEATABLE-READ

1 row in set(0.00 sec)

mysql＞SELECT*FROM parent WHERE id=1;

+----+

|id|

+----+

|1|

+----+

1 row in set(0.00 sec)

---

从这个小测试可以看出，在不同事务隔离级别下，6-8给出的存储过程测试所产生的结果也是不同的。在READ COMMITTED级别下，事务并不会得到较早时间的版本。此外，从某种角度上说，这也解释了锁升级（ACID中的I）的另一种表现形式。在第7章，我们还会继续讨论。

表 6-8　示例执行的过程

| 时间 | 会话 A | 会话 B |
| --- | --- | --- |
| 1 | BEGIN | |
| 2 | SELECT * FROM parent WHERE id = 1; | |
| 3 | | BEGIN |
| 4 | | UPDATE parent SET id=3 WHERE id = 1; |
| 5 | SELECT * FROM parent WHERE id = 1; | |
| 6 | | COMMIT; |
| 7 | SELECT * FROM parent WHERE id = 1; | |
| 8 | COMMIT | |

## 6.3.3 □□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□REPEATABLE READ□□□□□InnoDB□□□□□□SELECT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SELECT□□□□□□□InnoDB□□□□□□SELECT□□□□□□□□□□□□□□□□locking read□□□□□

□SELECT…FOR UPDATE

□SELECT…LOCK IN SHARE MODE

SELECT…FOR UPDATE□□□□□□□□□□□□□X□□□□□□□□□□□□□□□□□□□□□□□□□SELECT…LOCK IN SHARE MODE□□□□□□□□□□□□S□□□□□□□□□□□□□□□S□□□□□□□□X□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□SELECT…FOR UPDATE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SELECT…FOR UPDATE□SELECT…LOCK IN SHARE MODE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SELECT□□□□□□□□□□□BEGIN□START TRANSACTION□□□SET AUTOCOMMIT=0□

# 6.3.4　自增长与锁

自增长在数据库中是非常常见的一种属性，也是很多DBA或开发人员头疼的问题。在InnoDB存储引擎的内存结构中，对每个含有自增长值的表都有一个自增长计数器（auto-increment counter）。当对含有自增长的计数器的表进行插入操作时，这个计数器会被初始化，执行如下的语句来得到计数器的值：

---

```
SELECT MAX(auto_inc_col)FROM t FOR UPDATE;
```

---

插入操作会依据这个自增长的计数器值加1赋予自增长列。这个实现方式称做AUTO-INC Locking。这种锁其实是采用一种特殊的表锁机制，为了提高插入的性能，锁不是在一个事务完成后才释放，而是在完成对自增长值插入的SQL语句后立即释放。

虽然AUTO-INC Locking从一定程度上提高了并发插入的效率，但还是存在一些性能上的问题。首先，对于有自增长值的列的并发插入性能较差，事务必须等待前一个插入的完成（虽然不用等待事务的完成）。其次，对于INSERT…SELECT的大数据量的插入会影响插入的性能，因为另一个事务中的插入会被阻塞。

从MySQL 5.1.22版本开始，InnoDB存储引擎中提供了一种轻量级互斥量的自增长实现机制，这种机制大大提高了自增长值插入的性能。并且从该版本开始，InnoDB存储引擎提供了一个参数innodb_autoinc_lock_mode来控制自增长的模式，该参数的默认值为1。在继续讨论新的自增长实现方式之前，我们需要了解下自增长的分类，如表6-9所示。

## 表 6-9　插入类型

| 插入类型 | 说明 |
|---|---|
| insert-like | insert-like 指所有的插入语句，如 INSERT、REPLACE、INSERT…SELECT，REPLACE…SEECT、LOAD DATA 等 |
| simple inserts | simple inserts 指能在插入前就确定插入行数的语句。这些语句包括 INSERT、REPLACE 等。需要注意的是：simple inserts 不包含 INSERT …ON DUPLICATE KEY UPDATE 这类 SQL 语句 |
| bulk inserts | bulk inserts 指在插入前不能确定得到插入行数的语句，如 INSERT…SELECT，REPLACE…SELECT，LOAD DATA |
| mixed-mode inserts | mixed-mode inserts 指插入中有一部分的值是自增长的，有一部分是确定的。如 INSERT INTO t1 (c1,c2) VALUES (1,'a'), (NULL,'b'), (5,'c'), (NULL,'d'); 也可以是指 INSERT …ON DUPLICATE KEY UPDATE 这类 SQL 语句 |

参数 innodb_autoinc_lock_mode 的默认值为1，我们来具体分析每一种模式。参数 innodb_autoinc_lock_mode 以及这三种插入类型之间的关系，可以从0、1、2的取值详见表 6-10 所示。

表 6-10　参数 innodb_autoinc_lock_mode 的说明

| nnodb_autoinc_lock_mode | 说明 |
| --- | --- |
| 0 | 这是 MySQL5.1.22 版本之前自增长的实现方式，即通过表锁的 AUTO-INC Locking 方式。因为有了新的自增长实现方式，0 这个选项不应该是新版用户的首选项 |
| 1 | 这是该参数的默认值。对于"simple inserts"，该值会用互斥量（mutex）去对内存中的计数器进行累加的操作。对于"bulk inserts"，还是使用传统表锁的 AUTO-INC Locking 方式。在这种配置下，如果不考虑回滚操作，对于自增值列的增长还是连续的。并且在这种方式下，statement-based 方式的 replication 还是能很好地工作。需要注意的是，如果已经使用 AUTO-INC Locing 方式去产生自增长的值，而这时需要再进行"simple inserts"的操作时，还是需要等待 AUTO-INC Locking 的释放 |
| 2 | 在这个模式下，对于所有"INSERT-like"自增长值的产生都是通过互斥量，而不是 AUTO-INC Locking 的方式。显然，这是性能最高的方式。然而，这会带来一定的问题。因为并发插入的存在，在每次插入时，自增长的值可能不是连续的。此外，最重要的是，基于 Statement-Base Replication 会出现问题。因此，使用这个模式，任何时候都应该使用 row-base replication。这样才能保证最大的并发性能及 replication 主从数据的一致 |

此种方案的优点是利用了InnoDB的事务性特点以及利用MyISAM的高速。MyISAM对只读类的查询有很好的支持。在这种架构下，master选择InnoDB引擎，而slave采用MyISAM引擎。通过replication功能将数据同步到多台从服务器。

需要注意InnoDB表与其他存储引擎表在一些索引实现上的区别，如在建立的联合索引时，对于自增长列的，对于MySQL数据库类似如下的MyISAM存储引擎表，在建立联合自增长索引时是允许的，如：

---

mysql：CREATE TABLE t(

-：a INT AUTO_INCREMENT,

-：B INT,

-：KEY(b,a)

-：)ENGINE=InnoDB;

ERROR 1075(42000):Incorrect table definition;there can be only one auto column and it must be defined as a key

mysql：CREATE TABLE t(

-：a INT AUTO_INCREMENT,

-：B INT,

-：KEY(b,a)

-：)ENGINE=MyISAM;

Query OK,0 rows affected(0.01 sec)

---

# 6.3.5　外键和锁

外键主要用于引用完整性的约束检查。在InnoDB存储引擎中，对于一个外键列，如果没有显式地对这个列加索引，InnoDB存储引擎自动对其加一个索引，因为这样可以避免表锁——这比Oracle数据库做得好，Oracle数据库不会自动添加索引，用户必须自己手动添加，这也导致了Oracle数据库中可能产生死锁。

对于外键值的插入或更新，首先需要查询父表中的记录，即SELECT父表。但对于父表的SELECT操作，不是使用一致性非锁定读的方式，因为这样会发生数据不一致的问题，因此这时使用的是SELECT...LOCK IN SHARE MODE方式，即主动对父表加一个S锁。如果这时父表上已经加了X锁，子表上的操作会被阻塞，如表6-11所示。

表 6-11 外键测试用例

| 时　　间 | 会话 A | 会话 B |
| --- | --- | --- |
| 1 | BEGIN | |
| 2 | DELETE FROM parent WHERE id=3; | |
| 3 | | BEGIN |
| 4 | | INSERT INTO child SELECT 2,3<br>#第二列是外键，执行该句时被阻塞<br>(waiting) |

在上述的例子中，会话 B 会被阻塞。这里会话 B 被阻塞。这里会话 B 被阻塞。这里会话 B 被阻塞。这里会话 B 通过 COMMIT、ROLLBACK 来解锁。会话 B 希望插入的外键值id为3，而会话 A 已经对父表的X锁，因此会话 B 希望插入的外键值id为3所以被

加的S锁（由INSERT操作触发，用于唯一性判断冲突检测），因此两个事务的锁存在冲突。Session B试图更新id=3的记录，一直处于等待状态。如果在会话A中此时执行删除或更新id为3的记录，那么就可能触发死锁。下面我们通过查看INNODB_LOCKS视图来分析锁的情况。

mysql□SELECT*FROM information_schema.INNODB_LOCKS\G;

***************************1.row***************************

lock_id:7573B8:96:3:4

lock_trx_id:7573B8

lock_mode:S

lock_type:RECORD

lock_table:'mytest'.'parent'

lock_index:'PRIMARY'

lock_space:96

lock_page:3

lock_rec:4

lock_data:3

***************************2.row***************************

lock_id:7573B3:96:3:4

lock_trx_id:7573B3

lock_mode:X

lock_type:RECORD

lock_table:'mytest'.'parent'

lock_index:'PRIMARY'

lock_space:96

lock_page:3

lock_rec:4

lock_data:3

2 rows in set(0.00 sec)

# 6.4 锁的算法

## 6.4.1 行锁的3种算法

InnoDB存储引擎3种行锁的算法，其分别是：

□Record Lock：单个行记录上的锁

□Gap Lock：间隙锁，锁定一个范围，但不包含记录本身

□Next-Key Lock:Gap Lock+Record Lock，锁定一个范围，并且锁定记录本身

Record Lock总是会去锁住索引记录，如果InnoDB存储引擎表在建立的时候没有设置任何一个索引，那么这时InnoDB存储引擎会使用隐式的主键来进行锁定。

Next-Key Lock是结合了Gap Lock和Record Lock的一种锁定算法，在Next-Key Lock算法下，InnoDB对于行的查询都是采用这种锁定算法。例如一个索引有10，11，13和20这四个值，那么该索引可能被Next-Key Locking的区间为：

---

(-∞,10]

(10,11]

(11,13]

(13，20]

(20,+∞)

---

采用Next-Key Lock的锁定技术称为Next-Key Locking。其设计的目的是为了解决Phantom Problem，这将在下一小节中介绍。而利用这种锁定技术，锁定的不是单个值，而是一个范围，是谓词锁（predict lock）的一种改进。除此之外还有previous-key locking技术。同样上述的索引10，11，13和20，若采用previous-key locking技术，那么可锁定的区间为：

---

(-∞,10)

[10,11)

[11,13)

[13，20)

[20,+∞)

---

## 对于T1，根据next-key locking原则，锁住的区间：

---

(10,11]，(11，13]

---

## 由于隔离级是12，锁区间按照该原则缩小：

---

(10,11]，(11,12]，(12，13]

---

## 这里为什么要指出隔离级别呢？因为InnoDB默认实现的Next-Key Lock算法会将其降为Record Lock，即仅锁住索引本身，而非范围。下面是具体的例子。首先是创建表t：

---

```
DROP TABLE IF EXISTS t;

CREATE TABLE t(a INT PRIMARY KEY);

INSERT INTO t SELECT 1;

INSERT INTO t SELECT 2;

INSERT INTO t SELECT 5;
```

---

## 然后是如表6-12所示SQL语句：

## 表 6-12　唯一索引的锁定示例

| 时　间 | 会话 A | 会话 B |
|---|---|---|
| 1 | BEGIN; | |
| 2 | SELECT * FROM t<br>WHERE a =5 FOR UPDATE; | |
| 3 | | BEGIN; |
| 4 | | INSERT INTO t SELECT 4; |
| 5 | | COMMIT;<br>＃成功，不需要等待 |
| 6 | COMMIT | |

表t在时间1、2、5中并没有发生任何改变。会话A中通过对a＝5加上X锁，由于列a是主键且唯一，因此只锁定5这条记录，其范围(2，5)并没有被锁定。因此会话B在时间4可以进行插入操作，不会发生任何阻塞，即此时使用的是Record Lock，而不是之前的Next-Key Lock，这就提高了应用的并发性。

从上述的例子中看似Next-Key Lock的降级为Record Lock可以提高应用的并发性，但是如果在实际应用中使用不好，则可能导致另一个问题的产生，如建立如下的表z：

---

```
CREATE TABLE z(a INT,b INT,PRIMARY KEY(a),KEY(b));

INSERT INTO z SELECT 1,1;
```

```
INSERT INTO z SELECT 3,1;

INSERT INTO z SELECT 5,3;

INSERT INTO z SELECT 7,6;

INSERT INTO z SELECT 10,8;
```

在z表上b列上执行如下会话A中所示的SQL语句：

```
SELECT*FROM z WHERE b=3 FOR UPDATE
```

很明显，这时的SQL语句通过索引列b进行查询，因此其使用传统的Next-Key Locking技术加锁，并且由于存在两个索引，需要分别进行锁定。对于聚集索引，其仅对列a等于5的索引加上Record Lock。而对于辅助索引，其加上的是Next-Key Lock，锁定的范围是(1，3)，特别需要注意的是，InnoDB存储引擎还会对辅助索引下一个键值加上gap lock，即还有一个辅助索引范围为(3，6)的锁。因此若在新会话B中运行下面的SQL语句，都会被阻塞：

```
SELECT*FROM z WHERE a=5 LOCK IN SHARE MODE;

INSERT INTO z SELECT 4,2;

INSERT INTO z SELECT 6,5;
```

第一个SQL语句不能执行，因为在会话A中执行的SQL语句已经对聚集索引中列a=5的值加上X锁，因此执行会被阻塞。第二个SQL语句，主键插入4，没有问题，但是插入的辅助索引值2在锁定的范围(1，3)中，因此执行同样会被阻塞。第三个SQL语句，插入的主键6没有被锁定，5也不在范围(1，3)之间。但插入的值5在另一个锁定的范围(3，6)中，故同样会被阻塞。下面的SQL语句，不会被阻塞，可以立即执行：

```
INSERT INTO z SELECT 8,6;

INSERT INTO z SELECT 2,0;

INSERT INTO z SELECT 6,7;
```

从上面的例子中可以看到，Gap Lock的作用是为了阻止多个事务将记录插入到同一范围内，而这会导致Phantom Problem问题的产生。例如在上面的例子中，会话A中用户已经锁定了b=3的记录。若此

因此，Gap Lock锁定（3，6），这样在b字段3两侧的间隙被封锁，A事务再执行这样
两条语句就不会产生影响，也就避免了Phantom Problem幻象问题。

可以通过以下两种方式来关闭Gap Lock：

□ 将事务的隔离级别设置为READ COMMITTED

□ 将参数innodb_locks_unsafe_for_binlog设置为1

在上述设置下，除了外键约束和唯一性检查依然需要Gap Lock，其余情况仅使用Record Lock
进行锁定。但需要牢记的是，上述设置破坏了事务的隔离性，并且对于replication，可能会导致主
从数据的不一致。此外，从性能上来看，READ COMMITTED也不会优于默认的隔离级别READ
REPEATABLE。

在InnoDB存储引擎中，对于INSERT的操作，其会检查插入记录的下一条记录是否被锁定，
若已经被锁定，则不允许查询。对于上面的例子，会话A已经锁定了表z中b=3的记录，即已经锁定了(1，3)之间的间隙，
这时若在其他会话中插入如下记录都会导致阻塞：

```
INSERT INTO z SELECT 2,2;
```

因为在区间（1，3）之间插入2记录，故将被锁定。但是如果插入3记录，则可以立即插入成功，
不会产生锁等待：

```
INSERT INTO z SELECT 2,0;
```

从上面的例子中可以看到，Gap Lock的作用是为了阻止多个事务将记录插入到同一范围内，
而这会导致Phantom Problem问题的产生。例如在上面的例子中，会话A中用户已经锁定了b=3的
记录。若此时没有Gap Lock锁定（3，6），那么用户可以插入索引b列为3的记录，这会导致会话A
中的用户再次执行同样查询时会返回不同的记录，导致Phantom Problem问题的产生。

用户可以通过以下方式来显示关闭Gap Lock：将事务的隔离级别设置为READ COMMITTED，或将参数innodb_locks_unsafe_for_binlog设置为1。

对于唯一键值的锁定，Next-Key Lock降级为Record Lock仅存在于查询所有的唯一索引列。
若唯一索引由多个列组成，而查询仅是查找多个唯一索引列中的其中一个，那么查询其实是
range类型查询，而不是point类型查询，故InnoDB存储引擎依然使用Next-Key Lock进行
锁定。

## 6.4.2 　解决Phantom Problem

在默认的事务隔离级别下，即REPEATABLE READ下，InnoDB存储引擎采用Next-Key Locking机制来避免Phantom Problem（幻像问题）。这点可能不同于其他的数据库，如Oracle数据库，因为其可能需要在SERIALIZABLE的事务隔离级别下才能解决Phantom Problem。

Phantom Problem是指在同一事务下，连续执行两次同样的SQL语句可能导致不同的结果，第二次的SQL语句可能会返回之前不存在的行。下面通过一个例子来说明：假设表t有如下记录，t有1、2、5行，在默认的隔离级别下，即T1执行如下SQL语句：

```
SELECT*FROM t WHERE a≥2 FOR UPDATE;
```

我们注意到事务T1还没有进行提交操作，表中有一个记录5大于查询的条件。这时假设另一个事务T2插入了4这条记录，并且提交了该事务。回到事务T1执行上述的SQL语句又会获得4和5这两个记录，这样在同一个事务中两次执行同样的查询，却会得到不同的结果，这种现象称为幻像，如图6-13所示。

表 6-13 Phantom Problem 的演示

| 时　间 | 会话 A | 会话 B |
|---|---|---|
| 1 | SET SESSION<br><br>tx_isolation='READ-OMMITTED'; | |
| 2 | BEGIN; | |
| 3 | SELECT * FROM t<br><br>WHERE a > 2 FOR UPDATE;<br><br>********** 1. row *********<br><br>a: 4 | |

| 时　间 | 会话A | 会话B |
|---|---|---|
| 4 | | BEGIN; |
| 5 | | INSERT INTO t SELECT 4; |
| 6 | | COMMIT; |
| 7 | SELECT * FROM t<br><br>WHERE a > 2 FOR UPDATE;<br><br>********* 1.row ********<br><br>a: 4<br><br>********* 2.row ********<br><br>a: 5 | |

InnoDB存储引擎采用Next-Key Locking的算法避免Phantom Problem。对于上述的SQL语句SELECT*FROM t WHERE a﹥2 FOR UPDATE，其锁住的不是5这单个值，而是对（2，+∞）这个范围加了X锁。因此任何对于这个范围的插入都是不被允许的，从而避免Phantom Problem。

InnoDB存储引擎默认的事务隔离级别为REPEATABLE READ，在该隔离级别下，其采用Next-Key Locking的方式来加锁。而在隔离级别READ COMMITTED下，其仅采

Record Lock锁定的是单个行记录，但是A如果事务的隔离级别为READ COMMITTED，

在某些情况下，InnoDB会使用到Next-Key Locking，以下是使用的一种典型场景

---

```
SELECT*FROM table WHERE col=xxx LOCK IN SHARE MODE；

If not found any row:

#unique for insert value

INSERT INTO table VALUES(...);
```

---

上述这种情况，很明显是多个事务并发地对SLock锁定的资源进行插入操作，这样就会产生死锁，因此不论隔离级别是何种，只要执行了SELECT... LOCK IN SHARE MODE，就会对不存在的值加上一个间隙锁，因此多个事务之间就不会并发地对同一资源进行插入操作，也就避免了死锁，如图6-14所示。

表 6-14  通过 Next-Key Locking 实现应用程序的唯一性检查

| 时  间 | 会话 A | 会话 B |
|---|---|---|
| 1 | BEGIN | |
| 2 | mysql>SELECT * FROM z<br>WHERE b=4<br>LOCK IN SHARE MODE; | |
| 3 | | mysql>SELECT * FROM z<br>WHERE b=4<br>LOCK IN SHARE MODE; |
| 4 | mysql>INSERT INTO z SELECT 4,4;<br>＃阻塞 | |
| 5 | | mysql>INSERT INTO z<br>SELECT4,4;<br>ERROR 1213 (40001):Deadlock found when<br>trying to get lock;try restarting transaction<br>＃抛出死锁异常 |
| 6 | ＃INSERT 插入成功 | |

## 6.5 隔离性

数据库系统中将并发执行的多个事务之间相互隔离的性质称为隔离性。有些情况下为了追求系统的并发性，会相应降低事务之间的隔离程度，由此引入了不同的隔离级别。

## 6.5.1 问题

脏读（又称为Dirty Read）。对于两个事务而言，当其中一个事务读到了另一个事务尚未提交的数据时就会产生脏读问题。比如说有两个事务，一个事务正在对某条数据进行修改操作，但是这时候尚未提交，而另一个事务却读取到了这条修改了但尚未提交的数据，由此产生脏读问题。脏读问题主要发生在事务尚未commit时。

不可重复读。对于两个事务而言，当其中一个事务对数据进行多次读取的时候，而在读取的间隙中另一个事务对数据进行了修改，由此导致前一个事务两次读取到的结果不一致，从而产生不可重复读问题。

幻读。与不可重复读类似，幻读问题同样是一个事务进行多次读取操作，在读取的间隙中另一个事务对数据进行了增删操作。

下面列出了不同的隔离级别所对应解决的不同问题，读者通过隔离级别可知道它们所要解决的问题，如表6-15所示（一般默认为可重复读）。

表 6-15 脏读的示例

| Time | 会话 A | 会话 B |
|------|-------|-------|
| 1 | SET<br>@@tx_isolation='read-ncommitted'; | |
| 2 | | SET<br>@@tx_isolation='read-ncommitted'; |
| 3 | | BEGIN; |
| 4 | | mysql> SELECT * FROM t\G;<br>********* 1. row ***********<br>a: 1<br>1 row in set (0.00 sec) |
| 5 | INSERT INTO t SELECT 2; | |
| 6 | | mysql> SELECT * FROM t\G;<br>********* 1. row ***********<br>a: 1<br>********* 2. row ***********<br>a: 2<br>2 row in set (0.00 sec) |

表t中）。在小节6.4.1的示例中就有这样的一种现象，在当前读一致性下可能出现问题。当隔离级别为REPEATABLE READ时，READ UNCOMMITTED对应的问题中，事务A所读取到的数据，可能会因为事务B中提交的SELECT操作而发生变化，如果在2个事务中，事务A使用当前读一致性，就可能导致上述问题的发生。

对于数据库系统而言，可能提供了不同的隔离级别供用户选择，不同的数据库系统可能会提供不同的默认隔离级别。在实际使用中，READ UNCOMMITTED很少被使用，大多数数据库的默认隔离级别为READ COMMITTED，InnoDB存储引擎的默认隔离级别为READ REPEATABLE，Microsoft SQL Server默认为READ COMMITTED，Oracle数据库默认为READ COMMITTED。

在实际应用中，也可能需要根据实际的业务需求来调整隔离级别，例如将隔离级别设置为READ UNCOMMITTED，在replication中，对于slave而言，如果将slave的隔离级别设置为比主库更低的级别。

## 6.5.2 　数据的处理

数据库的操作无非是对数据的处理，比较常见的对数据的处理是插入、删除和更新数据的操作，也就是所谓的 DML。当然，在网页中直接对数据进行处理是不大现实的，所以通常对数据的操作都是通过表单来进行的，也就是说通过对表单的处理来实现对数据库中数据的处理。

关于表单这一块主要涉及两个文件，一个是表单的页面，也就是供用户输入数据的界面；另一个是对提交数据进行处理的页面，也就是具体进行数据库操作的页面。如图6-16所示。

表 6-16 不可重复读的示例

| Time | 会话 A | 会话 B |
|---|---|---|
| 1 | SET@@tx_isolation='read-committed'; | |
| 2 | | SET @@tx_isolation='read-committed'; |
| 3 | BEGIN | BEGIN |
| 4 | mysql>SELECT * FROM t ;<br><br>********* 1. row ***********<br><br>a: 1<br><br>1 row in set (0.00 sec) | |
| 5 | | INSERT INTO t SELECT 2; |
| 6 | | COMMIT; |
| 7 | mysql>SELECT * FROM t ;<br><br>********* 1. row ***********<br><br>a: 1<br><br>********* 1. row ***********<br><br>a: 2<br><br>2 row in set (0.00 sec) | |

事务A第一次读取的数据和提交的数据为1，而第二次事务B又读取的同一行的数据为2。会导致这种问题的原因是事务A在两次读取相同行数据1的过程中，被其他事务B修改了这行数据。因此事务A两次读取相同行数据1和2出现了不一致的现象。这种隔离级别下，事务A读到了B提交的数据，隔离级别为READ COMMITTED。

这种隔离级别下的数据库系统，读取数据并不需要加锁，只有修改数据才需要加锁，因此很多数据库如Oracle、Microsoft SQL Server的默认事务隔离级别都是READ COMMITTED，但这种隔离级别会导致幻读现象。

由InnoDB存储引擎通过使用Next-Key Lock算法避免了幻读现象，所以MySQL的默认隔离级别就能避免幻读现象（Phantom Problem）的产生。Next-Key Lock是行锁与间隙锁的结合。当查询的索引含有唯一属性时会对gap间隙加锁，防止其他事务插入新的记录，这样当前事务再次读取数据时就不会出现幻读现象。所以InnoDB存储引擎的默认隔离级别是READ REPEATABLE，采用Next-Key Lock锁的算法可以避免幻读现象。

# 6.5.3　更新丢失

当两个或多个事务选择同一行，然后基于最初选定的值更新该行时，会发生更新丢失的问题，每个事务都不知道其他事务的存在。

1、事务T1读取行r，得到v1，此时T1未提交。

2、与此同时事务T2读取行r，得到v2，然后T2先提交。

3、此时T1再更新。

4、然后T2提交。

数据库管理系统为了实现并发性控制，最基本的手段就是加锁，通过锁，可以使得数据库事务在某个隔离级别下安全地运行。READ UNCOMMITTED事务隔离级别下，它并未对DML操作语句加锁，那么它的并发性控制是如何解决的呢？在步骤2，事务提交后T2再次读取同一行r，最后的更新会覆盖其他事务所做的更新，即被T1覆盖。

一个更为常见的更新丢失情况是更新丢失的另一种情况，即第二类型丢失更新。它是两个事务同时读取同一行数据，然后其中一个对它进行修改提交，而另一个也进行修改提交。这就会造成第一个事务所做的修改丢失。例如：

1、事务T1读取某商品的库存为某个值，此时用户为User1。

2、事务T2读取某商品的库存为某个值，此时用户为User2。

3、User1修改这个库存值，提交事务完成。

4、User2修改这个库存值，提交事务完成。

这样就出现了一种情况：User1所做的更新被"覆盖"，库存更新的值被另一个"覆盖"。比如最初的库存值为10 000，经过多个事务操作后，库存更新为9000，然而事务之间彼此没有进行协调，会出现这样一种情况，库存更新为1。为了避免这个问题，最常见的方案就是把库存更新为9999，而不是直接更新为9000。这个例子说明，在并发处理的过程中，多个事务操作同一条数据时，最终更新的值为9000，这个数据会造成很大的损失。如果你在网上购物时使用了USB Key进行加密，而在购买的过程中，USB Key出现了异常，那么你的付款就会出现问题，甚至会造成你多付款或者付款金额有误的情况。

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□1、□□□□□□□□□□□□□□□□□□□□□X□□□□□□□□□；2、□□□□□□□□□□□□□□□□□□□□□□□□□□□□X□□□□□□□□□□□□□□2□□□□□□□□□□□□1□□；（3）□□□□□□□□□□□；（4）□□（6-17□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

表 6-17　丢失更新问题的处理方法

| Time | 会话 A | 会话 B |
|---|---|---|
| 1 | BEGIN; | |
| 2 | SELECT cash into @cash<br>FROM account<br>WHERE user = pUser FOR UPDATE; | |
| 3 | | SELECT cash into @cash<br>FROM account<br>WHERE user = pUser FOR UPDATE;<br>＃等待 |
| | …… | …… |
| m | UPDATE account<br>SET cash=@cash-9000<br>WHERE user=pUser | |
| m+1 | COMMIT | |
| m+2 | | UPDATE account SET cash=@cash-1<br>WHERE user=pUser; |
| m+3 | | COMMIT |

□□□□□□□□□□□□□□□□□□□□□□□□UPDATE□□□□□□□□□□□SELECT...FOR UPDATE□□□□□□□□□□□□□UPDATE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□UPDATE□□□□□□□SELECT□UPDATE□□□□□□□□□□□□□□□SQL□□□□

□□□□□□□□□□□□□□□□SELECT□INSERT□UPDATE□DELETE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 6.6  □□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□InnoDB□□□□□□□□□innodb_lock_wait_timeout□□□□□□□□□□□□□□□□□□50□□□innodb_rollback_on_timeout□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□OFF□□□□□□□□□□□□innodb_lock_wait_timeout□□□□□□□□□□MySQL□□□□□□□□□□□□□□□

---

```
mysql□SET@@innodb_lock_wait_timeout=60;

Query OK,0 rows affected(0.00 sec)
```

---

## □innodb_rollback_on_timeout□□□□□□□□□□□□□□□□□□□□□□□□

---

```
mysql□SET@@innodb_rollback_on_timeout=on;

ERROR 1238(HY000):Variable'innodb_rollback_on_timeout'is a read only variable
```

---

## □□□□□□□MySQL□□□□□□□□□□1205□□□□□□□□□

---

```
mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□SELECT*FROM t WHERE a=1 FORUPDATE;

ERROR 1205(HY000):Lock wait timeout exceeded;try restarting transaction
```

---

## □□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

```
#□□A

mysql□SELECT*FROM t;
```

```
+---+

|a|

+---+

|1|

|2|

|4|

+---+

3 rows in set(0.00 sec)

mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□SELECT*FROM t WHERE a□4 FOR UPDATE;

+---+

|a|

+---+

|1|

|2|

+---+

2 rows in set(0.00 sec)
```

---

# 此时，A事务在上面使用了Next-Key Lock锁定了小于等于4的记录，以及记录4到正无穷之间的间隙锁。此时在B事务中执行下面的

---

```
#事务B

mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO t SELECT 5;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERTINTO t SELECT 3;

ERROR 1205(HY000):Lock wait timeout exceeded;try restarting transaction
```

若这时在会话事务B中删除了主键为5的这条记录，那么将会产生一个3的记录，会话A的Next-Key Lock尚没有释放，所以在会话A中再运行一次同样的一个当前的一个读的SELECT，可以看到主键5的这条记录已经不在了。

```
mysql＞SELECT*FROM t;

+---+
|a|

+---+
|1|

|2|

|4|

|5|

|8|

+---+

5 rows in set(0.00 sec)
```

这时再通过会话事务B试图删除这条记录，再来做一个COMMIT，不管是提交还是ROLLBACK，因为我们用的是快照读，所以不管是通过COMMIT还是ROLLBACK，其实对事务都是没有影响的。

# 6.7 □□

## 6.7.1 □□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□innodb_lock_wait_timeout□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□FIFO□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□undo log□□□□□□FIFO□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□wait-for graph□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□wait-for graph□□□□□□□□□□□□□□□□□□

□□□□□□□□□

□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□wait-for graph□□□□□□□□□□□□□□□□□□□□□□T1□□T2□□□□□□□

□□□□T1□□□□T2□□□□□□□

□□□□T1□□□□T2□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□T1□□□□□T2□□□□

□□□□□□□□□□□□□□□□□□□□□□□6-5□□□□

图 6-5 事务等待和锁申请表结构

在Transaction Wait Lists中，系统事务有4个，为t1、t2、t3、t4，若用wait-for graph表示这4个事务之间的关系，t2对row1加了x锁，此时t1对row2加了s锁，因此t1事务在等待t2，row1锁资源，在wait-for graph中表现为事务t1等待事务t2。由于t2事务在等待t1、t4事务对row2加的锁资源，因此t2等待t1、t4。因此，事务t3等待t1、t2、t4事务。相应的wait-for graph如6-6所示。

图 6-6 wait-for graph

在图6-6中，事务之间的①t1、t2？？？？？？？？？？？？？？？？？？？？？？wait-for graph？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？InnoDB？？？？？？？？？？undo？？？？？？？？？

wait-for graph？？？？？？？？？？？？？？？？？？？？？InnoDB1.2？？？？？？？？？？？？？？？？？？？1.2？？？？？？？wait-for graph？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？？InnoDB？？？？？？？？？

# 6.7.2 死锁概率

死锁是并发系统中最难以处理的问题之一，下面我们从理论上来分析在封锁协议下（特别是2阶段封锁协议），事务发生死锁的概率。只有极少数的数据库系统相关书籍谈及此话题。

假设每个事务平均获得n+1个锁，也就是说一个事务会更新n+1条记录，那么这个事务持有锁的平均数量为r+1。假设数据库中共有R条记录，并假设每条记录只被少量事务更新，且每个事务访问的记录互不关联。那么平均每个事务持有锁的数量为nr（即R之中的某部分），那么整个系统持有锁的总数量为：

当一个事务申请封锁一条记录时，需要等待的概率为已被其他事务封锁的记录占总记录的比例。一个事务持有锁的平均数量为：

$$(1+2+3+\cdots+r) / (r+1) \approx r/2$$

而数据库中一共有R条记录，所以一条记录被某个事务封锁的概率为1/R，那么申请封锁一条记录时需要等待的概率PW为：

$$PW=nr/2R$$

每个事务有r个锁，所以一个事务需要等待的概率PW(T)为：

$$PW(T)=1-(1-PW)^r \approx r*PW \approx \frac{nr^2}{2R}$$

一个事务等待另一个事务，而另一个事务又等待这个事务，这就是2阶段死锁，所以一个事务发生死锁的概率为：

$$一个事务发生死锁的概率 \approx \frac{PW(T)^2}{n} \approx \frac{nr^4}{4R^2}$$

从这个理论公式可以看出，发生2阶段死锁的概率是很低的，而且和事务的大小成正比，和数据库中记录数量成反比。

$$系统中任何一个事务发生死锁的概率 \approx \frac{n^2 r^4}{4R^2}$$

事务很少涉及其他事务，因此当n和R都较大或者r较小时，系统中事务发生死锁的可能性很小。上式中各项的含义如下：

☐ 系统中并发事务的数目n。事务数目越多，越容易发生死锁。

☐ 每个事务申请锁的步骤数r。步骤数越多，越容易发生死锁。

☐ 系统中数据项的数目R。数据项越多，越不容易发生死锁。

## 6.7.3 特征的表示

可用特征之间关系的内在结构进行分类信息的描述，并进而找到一种有利于分类的特征表示方式。例如，如图6-18所示，若根据两个特征之间的关系，由A推出B、B推出A，则可以将两个特征表示为AB-BA的形式。

表 6-18　死锁用例 1

| 时间 | 会话 A | 会话 B |
|---|---|---|
| 1 | BEGIN; | |
| 2 | mysql>SELECT * FROM t<br>WHERE a = 1 FOR UPDATE;<br>********* 1. row ***********<br>a: 1<br>1 row in set (0.00 sec) | BEGIN |
| 3 | | mysql>SELECT * FROM t<br>WHERE a = 2 FOR UPDATE;<br>********* 1. row ***********<br>a: 2<br>1 row in set (0.00 sec) |
| 4 | mysql>SELECT * FROM t WHERE a = 2<br>FOR UPDATE;<br>＃等待 | |
| 5 | | mysql>SELECT * FROM t WHERE a = 1<br>FOR UPDATE;<br>ERROR 1213 (40001): Deadlock found when<br>trying to get lock; try restarting transaction |

在这个数据库中，表B中存在一个外键1213，但是数据库中不存在这个外键的关联值，因此A表B表之间不会存在完整性约束。在InnoDB存储引擎中，外键的约束依然存在，因此这里表B中存在外键，A表中不存在2这个值。在这种情况下，B表中的列还是会执行A表中不存在的值，因此在6.6这个版本之后，InnoDB存储引擎在外键的约束下，如果父表中不存在，则InnoDB存储引擎会抛出异常，外键的约束失败，并返回错误代码1213，表示外键约束失败的错误信息。

Oracle数据库有一个成熟的解决方案，而在此之前，InnoDB存储引擎的外键约束是在存储引擎级别实现的，而不是在MySQL数据库级别上实现的。

---

mysql□CREATE TABLE p(

-□aINT,

-□PRIMARY KEY(a)

-□)ENGINE=InnoDB;

Query OK,0 rows affected(0.00 sec)

mysql□CREATE TABLE c(

-□bINT,

-□FOREIGH KEY(b)REFERENCES p(a)

-□)ENGINE=InnoDB;

Query OK,0 rows affected(0.00 sec)

mysql□SHOW INDEX FROM c\G;

***************************1.row***************************

Table:c

Non_unique:1

Key_name:b

Seq_in_index:1

Column_name:b

Collation:A

Cardinality:0

Sub_part:NULL

Packed:NULL

Null:YES

Index_type:BTREE

Comment:

1 row in set(0.00 sec)

mysql＞DROP INDEX b ON c;

ERROR 1553(HY000):Cannot drop index'b':needed in a foreign key constraint

---

通过上面的例子可以看到，想要删除外键列，其实在InnoDB存储引擎中也是一种索引。通过命令可以看到索引b仍然存在，而且并不能删除。

在某些情况下，与其说常常会设置一个外键来增加数据的关联性，一个列可以被X锁阻塞。可以通过阻塞由于S锁产生的操作，以提高并发性。以下我们来通过一个例子来了解表t上的阻塞情况。

---

```
CREATE TABLE t(

a INT PRIMARY KEY

)ENGINE=InnoDB;

INSERT INTO t VALUES(1),(2),(4),(5);
```

---

表t中存在列a，值为4，接下来我们观察图6-19中的事务情况。

## 表 6-19  死锁用例 2

| 时　间 | 会话 A | 会话 B |
|---|---|---|
| 1 | BEGIN; | |
| 2 | | BEGIN; |
| 3 | SELECT * FROM t<br><br>WHERE a = 4 FOR UPDATE; | |
| 4 | | SELECT * FROM t<br><br>WHERE a <= 4 LOCK IN SHARE MODE;<br><br>-- 等待 |
| 5 | INSERT INTO t VALUES(3);<br><br>-- ERROR 1213 (40001): Deadlock found when<br><br>trying to get lock; try restarting transaction | |
| 6 | | -- 事务获得锁，正常运行 |

可以看到，会话A中已经对记录4加上了X锁，而会话A中接着执行3的插入操作需要等待，因此会话B获得对记录4的S锁。接着会话执行插入操作，而此时插入操作由于1、2的关系需要等待，因此会话5插入的需要等待会话B中记录4占用的S锁资源。同时，会话对于插入操作3记录又等待会话中的记录进行。因此发生

InnoDB的崩溃恢复从宏观上看就是通过重做日志（undo log）进行的一个典型的AB-BA型死锁的标准实现。

# 6.8　□□□

□□□□Lock Escalation□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□1000□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

Microsoft SQL Server□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□Microsoft SQL Server 2005□□□□□□SQL Server□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□5000□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□40%□□□□□□□□□□□

□Microsoft SQL Server□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□3 000 000□□□□□□□□□□□□□□□□100□□□□□□□□□□□□300 000 000□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□X□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□10□□□□□□□□□□□□□□□□□□□□□□3GB□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□30□□□□□□□□□□□□□90MB□□□□□□□□□□□□□□□□□□□□□□□□□□

# 6.9　□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□

□□□□□□□□□□□□□□MySQL□□□InnoDB□□□□□□MyISAM□□□□□□Microsoft SQL Server□□□□□Oracle□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 第7章　事务

事务（Transaction）是数据库区别于文件系统的重要特性之一。在文件系统中，如果正在写文件，但是操作系统突然崩溃了，这个文件就很有可能被破坏。当然，有一些机制可以把文件恢复到某个一致性的状态。不过，如果需要保证两个文件同步，这些文件系统可能就显得无能为力了。例如，在需要更新两个文件时，更新完一个文件后，在更新完第二个文件之前系统重启了，就会有两个不同步的文件。

这正是数据库系统引入事务的主要目的：事务会把数据库从一种一致状态转换为另一种一致状态。在数据库提交工作时，可以确保要么所有修改都已经保存了，要么所有修改都不保存。

InnoDB存储引擎中的事务完全符合ACID的特性。ACID是以下4个词的缩写：

□ 原子性（atomicity）

□ 一致性（consistency）

□ 隔离性（isolation）

□ 持久性（durability）

第6章已经详细讨论了InnoDB存储引擎的隔离性，本章将关注事务的其他特性，以及使用事务时需要注意的一些地方。最后，给出一些作者关于事务的理解与心得。

# 7.1 事务处理

## 7.1.1 概述

事务处理在各种管理系统中SQL□□□□□□□□□□□□□□□，这些系统SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ACID□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ACID□□□□□□□□□□MySQL□NDB Cluster□□□□□□□□□□□□□□□□□□□□□□D□□□□□□□□□□□□□□□□Oracle□□□□□□□□□□□□□□□□□□□□READ COMMITTED□□□□□□I□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□READ REPEATABLE□□□□□□□□□□□□□□ACID□□□□□□□□□□□□□□□□ACID□□□□□□□□□□□□□□□□

A□Atomicity□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□C□□□□□□SQRT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQRT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ATM□□□□□□□□□□□□□□□□□□

1）用户ATM□□□□□□□□□□□□

2）□□□□□□□□□□□□□□□□□□□□□□□

3）□□□□□ATM□□□□□□□□□□□□□□□

4）□□□□□□□□□□□□□□□□□□□□□□□□

5）ATM□□□□□□

6）□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ATM□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□SQL□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□

C□consistency□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□——□□□□□□□□□

I□isolation□□□□□□□□□□□□□□□□□□□□□□□□□□□□□concurrency control□□
□□□□□serializability□□□□locking□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□granular lock□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□

D□durability□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□RAID□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□High Reliability□□□□□□□□□□□□□□High
Availability□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 7.1.2 　分类

根据事务处理过程的不同，将事务分为以下几类。

❑ 扁平事务（Flat Transactions）

❑ 带有保存点的扁平事务（Flat Transactions with Savepoints）

❑ 链事务（Chained Transactions）

❑ 嵌套事务（Nested Transactions）

❑ 分布式事务（Distributed Transactions）

扁平事务（Flat Transaction）是事务类型中最简单的一种，但在实际生产环境中，这可能是使用最为频繁的事务。在扁平事务中，所有操作都处于同一层次，其由BEGIN WORK开始，由COMMIT WORK或ROLLBACK WORK结束，其间的操作是原子的，要么都执行，要么都回滚。因此扁平事务是应用程序成为原子操作的基本组成模块。图7-1显示了扁平事务的有限状态机。

| | | |
|---|---|---|
| BEGIN WORK | BEGIN WORK | BEGIN WORK |
| Operation 1 | Operation 1 | Operation 1 |
| Operation 2 | Operation 2 | Operation 2 |
| …… | …… | …… |
| Operation K | （Error！！！） | |
| COMMIT WORK | ROLLBACK WORK | |
| | | 由于外界原因要回滚，如超时等 |
| 成功完成，约占所有事务的96% | 应用程序要求停止事务，约占所有事务的3% | |

强制终止事务，约占所有事务的1%

图 7-1  事务成功与失败示例

图7-1描述了数据库事务的三种结束情况。第一种是事务成功执行并提交，第二种是应用程序发现错误主动要求回滚，第三种是由于外界原因导致事务被强制终止并回滚。

数据库系统通过事务机制来保证数据的一致性和完整性。当事务成功完成时，所有操作的结果将被永久保存到数据库中；而当事务失败或被中止时，系统将撤销该事务所做的全部修改，使数据库恢复到事务开始之前的状态。

BEGIN WORK

S1：□□□□□□□□□□□

S2：□□□□□□□□□□□□□□□□□□

S3：□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□S3□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□S1、S2、S3□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□S1、S2□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□Flat Transactions with Savepoint□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Savepoint□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SAVE WORK□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□7-2□□□□□□□□□□□□□□□□

| BEGIN WORK |
| 隐含SAVE WORK：1 |
| Action |
| Action |
| SAVE WORK：2 |
| Action |
| SAVE WORK：3 |
| Action |
| Action |
| Action |
| SAVE WORK：4 |
| Action |
| ROLLBACK WORK：2 |

保存点2覆盖的操作

| Action |
| Action |
| SAVE WORK：5 |
| Action |
| SAVE WORK：6 |
| Action |
| Action |
| SAVE WORK：7 |
| Action |
| Action |
| ROLLBACK WORK：7 |

保存点5覆盖的操作

| Action |
| Action |
| SAVE WORK：8 |
| Action |
| COMMIT WORK |

# 图 7-2 嵌套的事务应用实例

图7-2描述的是一个嵌套事务应用的实例。这里，ROLLBACK WORK表示事务回退，即事务中止的恢复过程。BEGIN WORK表示事务开始，后续内容中的数字表示事务嵌套的层次。ROLLBACK WORK:2表示该事务回退到嵌套层次为2的位置，而不是回退到最外层。ROLLBACK WORK:7表示回退。COMMIT WORK表示事务提交，即事务中所有操作成功完成后的确认过程。

在嵌套事务执行过程中，每个子事务都有可能发生错误。图7-2所示的应用实例中，当嵌套层次为2的子事务执行发生错误，回退到嵌套层次为3的位置并重新执行时，该事务最终成功执行到嵌套层次为5的位置。ROLLBACK操作保证了在执行发生错误时，整个系统能够恢复到一个一致的、正确的状态。

在嵌套事务中，ROLLBACK WORK:2表示将事务回退到嵌套层次为2的位置，此时系统撤销嵌套层次为2以内的所有操作结果，但保留外层事务的操作结果，这就是ROLLBACK WORK。

链式事务（Chained Transaction）是一种特殊的事务处理方式，在这种方式中，事务执行的中间结果可以分为易失（volatile）结果和持久（persistent）结果两种类型，系统在处理过程中需要对这两种结果进行不同的管理和保护。

在分布式数据库系统中，事务处理机制需要考虑多个节点之间的协调与同步问题，以保证整个分布式系统的数据一致性和事务的正确执行，这种协调机制的具体实现方式将在图7-3中进行详细的介绍和说明。

图 7-3 　 一个事务的触发器激活另一个事务，构成运动的事务

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□COMMIT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□Nested Transaction□□□□□□□□□□□□□□□□□□□□□□top-level transaction□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□subtransaction□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□7-4□□□□

顶层事务　　　　　　子事务　　　　　　子事务　　　　　　子事务

BEGIN WORK
　　　　　　　　　　　　　　　　　　BEGIN WORK
　　　　　　　　　　BEGIN WORK　　　　......
调用子事务　　　　　调用子事务　　　　　COMMIT WORK
　　　　　　　　　　调用子事务
　　　　　　　　　　COMMIT WORK
　　　　　　　　　　　　　　　　　　BEGIN WORK　　　　BEGIN WORK
　　　　　　　　　　　　　　　　　　......　　　　　　......
　　　　　　　　　　BEGIN WORK　　　　COM MIT WORK　　COMMIT WORK
调用子事务　　　　　......
　　　　　　　　　　COMMIT WORK

　　　　　　　　　　BEGIN WORK
调用子事务　　　　　调用子事务　　　　　BEGIN WORK
　　　　　　　　　　调用子事务　　　　　......
COMMIT WORK　　　　COMMIT WORK　　　　COMMIT WORK

图 7-4　□□□□□□□□□□

□□□□Moss□□□□□□□□□□□

1□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

2□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

3□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□predecessor□□□□□□□□parent□□□□□□□□□□□□□□□□□□child□□□

4□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

5□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□A、C、I□□□□□□□□D□□□□□

□Moss□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□7-5□□□□

Tk

Tk11

Tk1

Tk121

Tk12

Tk2

Tk31

Tk3

S1

S11

S12

S121

S2

S3

S31

Commit

time

# 图 7-5 分布式数据库系统的事务结构

在图7-5中我们可以看到各组件间的调用关系。在一个事务中，每个事务（如Tk3）都有多个子事务（如S2）。每个子事务负责对特定数据库服务器进行操作。

当执行分布式事务时，需要保证各个子事务之间的协调一致。如果其中一个子事务操作失败，整个事务都应该回滚，以保证数据的一致性。在分布式系统中，事务的协调通常由一个协调者来完成。协调者负责管理各个参与者P₁之间的协调。X、Y两个参与者是由顶层的协调者P₁₁管理的，而P₁则负责管理顶层事务的协调。当顶层协调者P₁₁将管理权交给Z之后，这种协调关系被反向继承（counter-inherited）给顶层的P₁，这样3个参与者X、Y、Z就由同一个协调者管理，即P₁₂负责管理整个事务的协调一致。

在这个过程中，需要保证各个子事务之间的协调一致。如果其中一个子事务操作失败，整个事务都应该回滚，以保证数据的一致性。

分布式事务（Distributed Transactions）是指操作涉及多个数据库服务器的事务。下面通过一个例子来说明分布式事务。

假设要通过一台ATM机从某账户转账到另一个账户，转账金额为10 000。这个转账操作涉及多个步骤，这台ATM要先从账户A中取款，然后存入账户B中。这里假设账户A和B在C银行的不同数据库服务器上。

1）账户A余额足够支付；

2）账户B的余额增加了，增加了10 000；

3）账户C的余额减少了，减少了10 000；

4）账户A的余额减少了，账户A的余额减少了。

在这个转账过程中，如果账户A的余额不足，则整个事务应该回滚，以保证数据的一致性。分布式事务需要保证事务的ACID特性。在上面的例子中，如果第2步或第3步操作失败，整个事务都应该回滚，以保证数据的一致性。

□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

## 7.2 重做日志

回顾第一章，第6节中我已经介绍了事务的隔离性由锁来实现，redo log（undo log）来保证事务的原子性和持久性。redo log用来保证事务的持久性，undo log用来帮助事务回滚及

虽然DBA们都知道undo与redo是关于事务日志的，redo与undo的作用都可以视为是一种恢复操作，redo恢复提交事务修改的页操作，而undo回滚行记录到某个特定版本。因此两者记录的内容不同，redo通常是物理日志，记录的是页的物理修改操作。undo是逻辑日志，根据每行记录进行记录。

## 7.2.1 redo

### 1.基本概念

重做日志用来实现事务的持久性，即事务ACID中的D。其由两部分组成：一是内存中的重做日志缓冲（redo log buffer），其是易失的；二是重做日志文件（redo log file），其是持久的。

InnoDB是事务的存储引擎，其通过Force Log at Commit机制实现事务的持久性，即当事务提交（COMMIT）时，必须先将该事务的所有日志写入到重做日志文件进行持久化，待事务的COMMIT操作完成才算完成。这里的日志是指重做日志，在InnoDB存储引擎中，由两部分组成，即redo log和undo log。redo log用来保证事务的持久性，undo log用来帮助事务回滚及MVCC的功能。redo log基本上都是顺序写的，在数据库运行时不需要对redo log的文件进行读取操作。而undo log是需要进行随机读写的。

为了确保每次日志都写入重做日志文件，在每次将重做日志缓冲写入重做日志文件后，InnoDB存储引擎都需要调用一次fsync操作。由于重做日志文件打开并没有使用O_DIRECT选项，因此重做日志缓冲先写入文件系统缓存。为了确保重做日志写入磁盘，必须进行一次fsync操作。由于fsync的效率取决于磁盘的性能，因此磁盘的性能决定了事务提交的性能，也就是数据库的性能。

InnoDB存储引擎允许用户手工设置非持久性的情况发生，以此提高数据库的性能。即当事务提交时，日志不写入重做日志文件，而是等待一个时间周期后再进行fsync操作。由于并非强制在事务提交时进行一次fsync操作，显然这可以显著提高数据库的性能。但是当数据库发生宕机时，由于部分日志未刷新到磁盘，因此会丢失最后一段时间的事务。

参数innodb_flush_log_at_trx_commit用来控制重做日志刷新到磁盘的策略。该参数的默认值为1，表示事务提交时必须调用一次fsync操作。还可以设置该参数的值为0和2。0表示事务提交时不进行写入重做日志操作，这个操作仅在master thread中完成，而在master thread中每1秒会

如果允许待提交的事务fsync操作延迟2秒就会显著提高数据库性能。但这会带来一个问题,即若数据库发生宕机,由于部分事务没有从重做日志缓冲刷新到磁盘,因此会导致丢失最后一段时间的事务。

下面演示事务的异步提交对MySQL数据库性能的影响,这里将使用前面已经介绍过的存储过程。首先根据如下代码创建测试表t1和存储过程p_load:

---

```
CREATE TABLE test_load(

a INT,

b CHAR(80)

)ENGINE=INNODB;

DELIMITER//

CREATE PROCEDURE p_load(count INT UNSIGNED)

BEGIN

DECLARE s INT UNSIGNED DEFAULT 1;

DECLARE c CHAR(80)DEFAULT REPEAT('a',80);

WHILE s<=count DO

INSERT INTO test_load SELECT NULL,c;

COMMIT;

SET s=s+1;

END WHILE;

END;

//

DELIMITER;
```

---

存储过程p_load的作用是将数据不断地插入表test_load中,并且每插入一条就进行一次显式的COMMIT操作。在默认的设置下,即参数innodb_flush_log_at_trx_commit为1的情况下,InnoDB存储引擎会将重做日志缓冲中的日志写入文件,并调用一次fsync操作。如果执行命令CALL p_load(500 000),则会向表中插入50万行的记录,并执行50万次的fsync操作。先来看在默认情况下执行上述存储过程需要的时间:

---

```
mysql□CALL p_load(500000);

Query OK,0 rows affected(1 min 53.11 sec)
```

---

在上述测试中，50万行数据的插入约需要2分钟。下面将演示禁用二进制日志后，插入的时间同样也是非常缓慢的。这里通过设置参数fsync的次数，即来模拟二进制日志的刷新，将参数innodb_flush_log_at_trx_commit设置为0，即来说明：

---

```
mysql□SHOW VARIABLES LIKE'innodb_flush_log_at_trx_commit'\G

***************************1.row***************************

Variable_name:innodb_flush_log_at_trx_commit

Value:0

1 row in set(0.00 sec)

mysql□CALL p_load(500000);

Query OK,0 rows affected(13.90 sec)
```

---

上述实验将参数innodb_flush_log_at_trx_commit设置为0，同样是50万行数据，插入的时间缩短为了13.90秒，差不多仅为之前的12%。而这其中的最大的区别就在于减少了对于重做日志的fsync操作。因此为了提高数据库插入的性能，可以参考7-1节，将参数innodb_flush_log_at_trx_commit设置为其他的值，这样存储过程p_load的这50万行的插入操作将变得非常快。

表 7-1　不同 innodb_flush_log_at_trx_commit 设置对于插入的速度影响

| innodb_flush_log_at_trx_commit | 执行所用时间 |
| --- | --- |
| 0 | 13.90 秒 |
| 1 | 1 分 53.11 秒 |
| 2 | 23.37 秒 |

□□□□□□□□□□□□□innodb_flush_log_at_trx_commit□0□2□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ACID□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□50□□□□□□□□□□□□□□□□COMMIT□□□□□□□□□□□□□□□□□□□□□□□□□□□□COMMIT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□MySQL□□□□□□□□□□□□□□□□binlog□□□□□□□□□POINT-IN-TIME□PIT□□□□□□□□□□□□□Replication□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□7-6□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

binlog

| T1 | T4 | T3 | T2 | T8 | T6 | T7 | T5 |
|----|----|----|----|----|----|----|----|

reod log

| T1 | T2 | T1 | *T2 | T3 | T1 | *T3 | *T1 |
|----|----|----|-----|----|----|-----|-----|

图 7-6 二进制日志与重做日志的写入顺序关系

如图7-6所示，由于重做日志是顺序写入的，其记录的是对于每个页的修改。因此InnoDB存储引擎在恢复时，需要对每一个重做日志记录进行解析，然后按照页的物理位置进行重做，这是与二进制日志所不一样的地方。图中的*T1、*T2、*T3表示已经执行完成的事务。

## 2.log block

在InnoDB存储引擎中，重做日志都是以512字节进行存储的。这意味着重做日志缓存、重做日志文件都是以块（block）的方式进行保存的，称之为重做日志块（redo log block），每块的大小为512字节。

若一个页中产生的重做日志数量大于512字节，那么需要分割为多个重做日志块进行存储。此外，由于重做日志块的大小和磁盘扇区大小一样，都是512字节，因此重做日志的写入可以保证原子性，不需要doublewrite技术。

重做日志块除了日志本身之外，还由日志块头（log block header）及日志块尾（log block tailer）两部分组成。重做日志头一共占用12字节，重做日志尾占用8字节，故每个重做日志块实际可以存储的大小为492字节（512-12-8）。图7-7显示了重做日志块缓存的结构。

Redo Log Buffer

| log block | log block | log block | log block | log block | log block | log block | log block | ...... | log block | log block | log block |

log block header → LOG_BLOCK_HDR_NO

Log Block
512 Bytes

LOG_BLOCK_HDR_DATA_LEN

LOG_BLOCK_FIRST_REC_GROUP

LOG_BLOCK_CHECKPOINT_NO

log block tailer → LOG_BLOCK_TRL_NO

## 图 7-7 重做日志缓存的结构

图7-7中的内存区域很好理解，每个区域的大小均为512字节。若把重做日志缓存中的内存区域划分为log block header（前部空间）、log body（中部空间）、log block tailer（后部）。

log block header由4部分组成（见表7-2所示）

表7-2 log block header

| 名　称 | 占用字节 |
| --- | --- |
| LOG_BLOCK_HDR_NO | 4 |
| LOG_BLOCK_HDR_DATA_LEN | 2 |
| LOG_BLOCK_FIRST_REC_GROUP | 2 |
| LOG_BLOCK_CHECKPOINT_NO | 4 |

log buffer是由log block组成，在内部log buffer就好像一个数组。因此LOG_BLOCK_HDR_NO用来标记这个数组中的位置。其是递增并且循环使用的，占用4个字节，但是由于第一位用来判断是否是flush bit，所以最大的值为2G。

LOG_BLOCK_HDR_DATA_LEN占用2字节，表示log block所占用的大小。当log block被写满时，该值为0x200，表示使用全部log block空间，即占用512字节。

LOG_BLOCK_FIRST_REC_GROUP占用2个字节，表示log block中第一个日志所在的偏移量。如果该值的大小和LOG_BLOCK_HDR_DATA_LEN相同，则表示当前log block不包含新的日志。如事务T1的重做日志1占用762字节，事务T2的重做日志占用100字节。由于每个log block实际只能保存492个字节，因此其在log buffer中的情况应如图7-8所示。

图 7-8 LOG_BLOCK_FIRST_REC_GROUP的应用

如图7-8所示，假设事务日志记录T1占用的字节数为792，因此需要占用两个log block来存储。在第二个log block中，LOG_BLOCK_FIRST_REC_GROUP为12，即log block中第一个事务日志记录的偏移量。而在第二个log block中，还存储了事务日志记录T1的第二部分和事务T2，因此第二个log block的偏移量为第一个log block的LOG_BLOCK_FIRST_REC_GROUP为282（270+12）。

LOG_BLOCK_CHECKPOINT_NO占用4字节，表示该log block最后被写入时的检查点第4个字节的值。

log block tailer只有1个部分，如表7-3所示，其值和LOG_BLOCK_HDR_NO相同，并在log_block_init中被初始化。

表7-3　log block tailer 部分

| 名　称 | 大小（字节） |
| --- | --- |
| LOG_BLOCK_TRL_NO | 4 |

## 3.log group

log group为重做日志组，其中有多个重做日志文件。虽然源码中已支持log group的镜像功能，但是在ha_innobase.cc文件中禁止了该功能。因此InnoDB存储引擎实际只有一个log group。

log group是一个逻辑上的概念，并没有一个实际存储的物理文件来表示log group信息。log group由多个重做日志文件组成，每个log group中的日志文件大小是相同的，且在InnoDB 1.2版本之前，重做日志文件的总大小要小于4GB（不能等于4GB）。从InnoDB 1.2版本开始重做日志文件总大小的限制提高为了512GB。InnoSQL版本的InnoDB存储引擎从1.1版本开始又放松了4GB的限制。

重做日志文件中存储的就是之前在log buffer中保存的log block，因此其也是根据块的方式进行物理存储的管理，每个块的大小与log block块相同，同样为512字节。在InnoDB存储引擎运行过程中，log buffer根据一定的规则将内存中的log block刷新到磁盘。这个规则具体是：

❑事务提交时。

❑当log buffer中有一半的内存空间已经被使用时。

❑log checkpoint时。

对于log block的写入追加（append）在redo log file的最后部分，当一个redo log file被写满时，会接着写入下一个redo log file，其使用方式为round-robin。

一个log block。由于redo log file由多个（在早期版本中都是默认为一组的）redo log file组成，所以每一组redo log file都和log buffer一样，由log block组成。每一组第一个文件的前2KB部分存放的不是redo log file的前2KB部分存放的log block，而是整个log group的基础信息。redo log file的前2KB部分的前4个512字节存放的信息依次如表7-4所示。

高

表 7-4　redo log file 前 2KB 部分的内容

| 名　称 | 大小（字节） |
| --- | --- |
| log file header | 512 |
| checkpoint1 | 512 |
| 空 | 512 |
| checkpoint2 | 512 |

每组第一个文件的这些信息仅在每个log group的第一个redo log file中存放。在每个log group的其他redo log file中，虽然也保留了这部分，但是不会实际进行使用，即其他redo log file这部分内容也没有存放实际的log block，因而填充的2KB全部为空。正因为如此，InnoDB存储引擎在启动时总是会将第一个log group的redo log file的大小相加7-9所示。

## Log Group 1

Redo Log File1

| Log File Header | CP1 | | CP2 | Log Block | Log Block | Log Block | Log Block | Log Block | Log Block | ...... | Log Block |
|---|---|---|---|---|---|---|---|---|---|---|---|

Redo Log File2

| | | | | Log Block | Log Block | Log Block | Log Block | Log Block | Log Block | ...... | Log Block |
|---|---|---|---|---|---|---|---|---|---|---|---|

## Log Group 2

Redo Log File1

| Log File Header | CP1 | | CP2 | Log Block | Log Block | Log Block | Log Block | Log Block | Log Block | ...... | Log Block |
|---|---|---|---|---|---|---|---|---|---|---|---|

Redo Log File2

| | | | | Log Block | Log Block | Log Block | Log Block | Log Block | Log Block | ...... | Log Block |
|---|---|---|---|---|---|---|---|---|---|---|---|

图 7-9 log group与redo log file之间的关系

在log filer header后部分记录了InnoDB存储引擎保存的checkpoint的值，其作用将在下一小节中进行描述。这样可以看出，每个redo log文件组只有第一个文件的checkpoint的信息。

4.重做日志格式

不同的数据库操作会有对应的重做日志格式。此外，由于InnoDB存储引擎的存储管理是基于页的，故其重做日志格式也是基于页的。虽然有着不同的重做日志格式，但是它们有着通用的头部格式，如图7-10所示。



| redo_log_type | space | page_no | redo log body |

图 7-10 重做日志格式

通用的头部格式由以下3部分组成：

❑redo_log_type：重做日志的类型。

❑space：表空间的ID。

❑page_no：页的偏移量。

之后redo log body的部分，根据重做日志类型的不同，会有不同的存储内容，例如，对于页上记录的插入和删除操作，分别对应如图7-11所示的格式：

## MLOG_REC_INSERT

| type | space | page_no | cur rec _offset | len & extra_info | into_bits | origin_ offset | mis_matc h_index | rec body |
|------|-------|---------|------------------|-------------------|-----------|-----------------|-------------------|----------|

## MLOG_REC_DELETE

| type | space | page_no | offset |
|------|-------|---------|--------|

图 7-11 重做日志条目的结构

在InnoDB1.2版本时，一共有51种重做日志类型。随着功能不断地增加，相信会加入越来越多的重做日志类型。

## 5.LSN

LSN是Log Sequence Number的缩写，其代表的是日志序列号。在InnoDB存储引擎中，LSN占用8字节，并且单调递增。LSN表示的含义有：

❑重做日志写入的总量

❑checkpoint的位置

❑页的版本

LSN不仅记录在重做日志中，还存在于每个页中。在每个页的头部，有一个值FIL_PAGE_LSN，记录了该页的LSN。在页中，LSN表示该页最后刷新时LSN的大小。因为重做日志记录的是每个页的日志，因此页中的LSN用来判断页是否需要进行恢复操作。例如，页P1的LSN为10 000，而数据库启动时，InnoDB检测到写入重做日志中的LSN为13 000，并且该事务已经提交，那么数据库需要进行恢复操作，将重做日志应用到P1页中。同样的，对于重做日志中LSN小于P1页的LSN，不需要进行重做，因为P1页中的LSN表示页已经被刷新到该位置。

LSN表示事务写入重做日志的字节的总量。例如当前重做日志的LSN为1 000，有一个事务T1写入了100字节的重做日志，那么LSN就变为了1100，若又有事务T2写入了200字节的重做日志，那么LSN就变为了1 300。可见LSN记录的是重做日志的总量，其单位为字节。

LSN还有一个作用，用来记录每个数据页的版本信息。每个数据页的FIL_PAGE_LSN记录了该页的LSN。这里的LSN，指当前页最后被修改时的LSN值。该值越大说明当前页修改越晚。那如果数据页P1的LSN为10 000，重做日志中保存的InnoDB刷新到磁盘的LSN为13 000，那么就意味着需要将重做日志应用到该页，并且将当前页P1的版本升级到最新的版本。LSN在P1页的LSN小于重做日志保存的值，说明P1页的版本LSN较旧，需要升级。

这里可以通过命令SHOW ENGINE INNODB STATUS查看LSN的情况：

---

mysql□SHOW ENGINE INNODB STATUS\G;

……

---

LOG

---

Log sequence number 11 3047174608

Log flushed up to 11 3047174608

Last checkpoint at 11 3047174608

0 pending log writes,0 pending chkp writes

142 log i/o's done,0.00 log i/o's/second

……

1 row in set(0.00 sec)

---

Log sequence number表示当前的LSN，Log flushed up to表示刷新到重做日志文件的LSN，Last checkpoint at表示刷新到磁盘的LSN。

我们这里的例子中，Log sequence number和Log flushed up to的值是相同的，这是因为当前系统负载较低，可能刚刚执行过这些操作，或者此时数据库基本上处于没有业务操作的空闲时间。下面的例子就可以看到这几个值是不同的。

---

mysql□show engine innodb status\G;

---

LOG

```
---

Log sequence number 203318213447

Log flushed up to 203318213326

Last checkpoint at 203252831194

1 pending log writes,0 pending chkp writes

103447 log i/o's done,7.00 log i/o's/second

……

1 row in set(0.00 sec)
```

从上面的例子可以看到Log sequence number、Log flushed up to、Last checkpoint at三个值相差并不大。

## 6.小结

InnoDB存储引擎的重做日志对于数据库来说至关重要，它们记录了对于数据库的各种操作，以防止数据库发生宕机而导致的数据丢失。本章详细地介绍了InnoDB存储引擎重做日志的实现，并介绍了与重做日志相关的各种细节，其中重点介绍了两次写的原理。

由于checkpoint表示已经刷新到磁盘页上的LSN，因此在恢复过程中仅需恢复checkpoint开始的日志部分。如图7-12所示，当数据库在checkpoint的LSN为10 000时发生宕机，恢复操作仅恢复LSN 10 000～13 000范围内的日志。

图 7-12 　未刷新的页

InnoDB存储引擎会自动判断哪些页有必要进行刷新操作，哪些页仍可以保存在缓冲池中再进行缓存。例如对于上述的INSERT操作，表结构的定义如下所示，此时数据库需要插入的数据为：

```
CREATE TABLE t(a INT,b INT,PRIMARY KEY(a),KEY(b));
```

## 我们的SQL是这样

```
INSERT INTO t SELECT 1,2;
```

新插入的这条记录是需要存储到磁盘上的，我们看它存储的位置是：

```
page(2,3),offset 32,value 1,2#□□□□
page(2,4),offset 64,value 2#□□□□
```

这个偏移量的变化，可能是由于之前的数据的B+树的split导致页面分裂迁移导致的。不管怎么说，这个偏移量的变化是不影响的。

$$f(f(x)) = f(x)$$

所以DBA如果学过函数知识，那么就能比较容易理解这个ROW模式的幂等性。这里也有个小细节，就是这个INSERT语句要具有幂等性，就必须保证这个插入的位置是确定的，否则这个INSERT语句就不能具有幂等性。

# 7.2.2  undo

## 1.□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□"□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□undo□□□□□□□□□□□□□□□InnoDB□□□□□□□□□redo□□□□□□□□□□undo□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ROLLBACK□□□□□□□□□□□□□□□□□□undo□□□□□□□□□□□□□□□□□□□

redo□□□□□□□□□□□□□□redo□□□undo□□□□□□□□□□□□□□□□□□segment□□□□□□□□□□undo□□undo segment□□undo□□□□□□□□□□□□□□□□py_innodb_page_info.py□□□□□□□□□□□□□□□undo□□□□□□□□□□□□□□□□□□□□□ibdata1□□□2222□undo□□□

---

```
[root@xen-server□]#python py_innodb_page_info.py/usr/local/mysql/data/ibdata1

Total number of page:46208:

Insert Buffer Free List:13093

Insert Buffer Bitmap:3

System Page:5

Transaction system Page:1

Freshly Allocated Page:4579

undo Log Page:2222

File Segment inode:6

B-tree Node:26296

File Space Header:1

□□□□□□:2
```

---

□□□□□□undo□□□□□□□□□□undo□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□——□□□□□□□□□□undo□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□

□□□□□□□□□□□□INSERT 10W□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□ROLLBACK□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□INSERT□InnoDB□□□□□
□□□□□□DELETE□□□□□□DELETE□InnoDB□□□□□□□□□□□□□INSERT□□□□□□
UPDATE□InnoDB□□□□□□□□□□□□□□□UPDATE□□□□□□□□□□□

□□□□□□□□undo□□□□□□□□□□MVCC□□□□InnoDB□□□□□□MVCC□□□□□□□□□□
undo□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□undo□□□□□□□□
□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□undo log□□□□redo log□□□□□undo log□□□□□□□□□
redo log□□□□□□□□□□undo log□□□□□□□□□□□□□

## 2.undo□□□□

InnoDB□□□□□□undo□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
InnoDB□□□□□□□rollback segment□□□□□□□□□□□□□1024□undo log
segment□□□□□□undo log segment□□□□□undo□□□□□□□□□□□□□□□□□□5
□□□□□0□5□□□□□□□rollback segment header□□□□□□□□□□□□□□□□□□□
FIL_PAGE_TYPE_SYS□

□InnoDB1.1□□□□□□□□□□1.1□□□□□□□□□□rollback segment□□□□□□□□□□□
□□□□□□□□□□1024□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□1.1□□□□□
InnoDB□□□□□□128□rollback segment□□□□□□□□□□□□□□□□□□□□□□□□□□
128*1024□

□□□InnoDB1.1□□□□□□128□rollback segment□□□□□□□rollback
segment□□□□□□□□□□□□□□□□□□InnoDB1.2□□□□□□□□□□□□□□□rollback
segment□□□□□□□□□□□□□□□□□□□□□□

□innodb_undo_directory

□innodb_undo_logs

□innodb_undo_tablespaces

innodb_undo_directory：用于rollback segment文件所在路径。这意味着rollback segment可以存放在共享表空间以外的位置，即可以设置为独立表空间。该路径默认为"."，表示当前InnoDB存储引擎的目录。

　　innodb_undo_logs：用于设置rollback segment的个数，默认值为128。在InnoDB1.2版本中，该参数用来替换之前版本的innodb_rollback_segments。

　　innodb_undo_tablespaces：用于设置构成rollback segment文件的数量，这样rollback segment可以较为平均地分布在多个文件中。设置该参数后，会在路径innodb_undo_directory看到undo为前缀的文件，该文件就代表rollback segment文件。7-13显示了设置3个rollback segment。

```
myspl> SHOW VARIABLES LIKE 'innodb_undo%';

+------------------------+-------+
| Variable_name          | Value |
+------------------------+-------+
| innodb_undo_directory  | .     |
| innodb_undo_logs       | 128   |
| innodb_undo_tablespaces| 3     |
+------------------------+-------+
3 rows in set (0.00 sec)


mysql> SHOW VARIABLES LIKE 'datadir';

+--------------+-------------------------------+
| Variable_name| Value                         |
+--------------+-------------------------------+
| datadir      | /Users/david/mysql_data/data/ |
+--------------+-------------------------------+
1 row in set (0.00 sec)


mysql> system ls -lh/Users/david/mysql_data/data/undo*

-rw-rw----  1 david  staff    10M 11 22 16:55/Users/david/mysql_data/data/undo001

-rw-rw----  1 david  staff    10M 11 22 16:51/Users/david/mysql_data/data/undo002

-rw-rw----  1 david  staff    10M 11 22 16:51/Users/david/mysql_data/data/undo003
```

## 图 7-13 有3个页的链表的rollback segment

对于事务进行过程中产生的undo log segment，页需要支持undo log的数据在该事务提交后不能立即删除。这是因为可能其他事务需要通过InnoDB的事务隔离等级来读取之前版本的行记录。

当undo log需要被删除时，会被放入purge列表中。

对于undo log在页内的存储分布，我们也需要特别地注意。

事务在进行过程中会产生大量的undo log，undo log可能需要被存储在多个连续的页中。这些页按照链表的方式组织，也就是undo log链表，而在undo log链表中，其每一个undo log，undo log都具有purge的操作属性。

这里需要再次强调的是，undo页会被其他事务进行重用。对于OLTP的应用环境，这种情况非常常见。在这样的环境下，假设某应用的删除和更新操作的TPS（transaction per second）为1000，为每一个事务产生的undo日志量为1000*60字节，那么一分钟会产生1GB。若没有purge操作，则需要存储20分钟内的所有undo数据。这也说明了在删除、更新频繁的应用中，InnoDB存储引擎中的undo使用量是非常巨大的。为了避免这个问题，undo log被设计成可以被重用的，当有多个undo页被存储时，则只会占用当前undo页的3/4空间。若一个undo页被使用的超过了该大小，则会另外分配一个新的undo页来使用，旧的undo页被重用，则旧的undo页与新的undo页链接在一起，并放入undo页链表中。若该undo页被放入purge列表中，则意味着该页可以被重用于后续的事务。

我们可以通过SHOW ENGINE INNODB STATUS来查看当前undo log的相关信息。

---

```
mysql：SHOW ENGINE INNODB STATUS\G;

****************************1.row****************************

......

------------

TRANSACTIONS

------------

Trx id counter 3000

Purge done for trx's n:o：2C03 undo n:o：0

History list length 12
```

```
LIST OF TRANSACTIONS FOR EACH SESSION:

---TRANSACTION 0,not started

MySQL thread id 1,OS thread handle 0x1500f1000,query id 4 localhost root

show engine innodb status

......
```

History list length□□□□undo log□□□□□□□□12□purge□□□□□□□□□□□□□□□undo log□□□□□□□□□□□□□□□□□□□□History list length□□□□□□□□0□

## 3.undo log□□

□InnoDB□□□□□□□undo log□□□□

□insert undo log

□update undo log

insert undo log□□□insert□□□□□□undo log□□□insert□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□undo log□□□□□□□□□□□□□□□□□purge□□□insert undo log□□□□□□7-14□□□□

insert undo log record

| next | |
|---|---|
| type_cmpl | |
| *undo no | |
| *table id | |
| *len1 | col1 |
| *len2 | col2 |
| …… | …… |
| *lenN | colN |
| start | |

n_unique_index

图 7-14　insert undo log示意图

图7-14描述了insert undo log的组成，其中*表示对存储的数据进行了压缩。insert undo log开始部分的前两字节next记录的是下一个undo log的位置，通过next的字节可以得到下一个undo log的指针。结尾部分的两字节记录的是undo log开始的位置。接下来的type_cmpl，记录的是undo的类型，对于insert undo log，其值总是11。undo_no记录事务的ID，table_id记录undo log所对应的表对象。这两个值之后的部分，表示记录列的主键信息。在进行rollback操作时，根据这些值可以定位到具体的记录，然后进行删除即可。

update undo log记录的是对delete和update操作产生的undo log。该undo log可能需要提供MVCC机制，因此不能在事务提交时就进行删除。提交时放入undo log链表，等待purge线程进行最后的删除。update undo log的结构如图7-15所示。

update undo log record

| | | |
|---|---|---|
| next | | |
| type_cmpl | | |
| *undo no | | |
| *table id | | |
| info_bits | | |
| *DATA_TRX_ID | | |
| *DATA_ROLL_PTR | | |
| *len1 | i_col1 | |
| *len2 | i_col2 | |
| …… | …… | |
| *lenN | i_colN | |

n_unique_index

update vector

| | | |
|---|---|---|
| n_update_field | | |
| *pos1 | *len1 | u_old_col1 |
| *pos2 | *len2 | u_old_col2 |
| …… | …… | …… |
| *posN | *lenN | u_old_colN |
| n_bytes_below | | |
| *pos | *len | col1 |
| *pos | *len | col2 |
| …… | …… | …… |
| *pos | *len | colN |
| start | | |

表 7-15 update undo log类型

update undo log记录的前面部分与insert undo log相似，它也有表示记录长度的字段next，start，undo_no，table_id。但是，与insert undo log只有一种类型type_cmpl不同，update undo log会有以下三种不同的记录类型。

□12 TRX_UNDO_UPD_EXIST_REC：更新non-delete-mark的记录

□13 TRX_UNDO_UPD_DEL_REC：将delete的记录标记为not delete

□14 TRX_UNDO_DEL_MARK_REC：将记录标记为delete

接着的部分为update_vector信息，update_vector表示update操作导致发生改变的列。每个字段表示为一个列。由于update undo log记录的是undo log，因此需要对列进行更新的旧值记录。

## 4.查看undo信息

Oracle和Microsoft SQL Server数据库都由内部的数据字典来观察当前undo的信息，而InnoDB存储引擎在这方面做得还不够，DBA只能通过原理来推断当前的状态。InnoSQL对information_schema进行了扩展，添加了两张数据字典表，这样用户可以非常方便和快捷地查看undo的信息。

首先创建的数据字典表为INNODB_TRX_ROLLBACK_SEGMENT，该表用来查看当前数据库中回滚段rollback segment的信息，如表7-16所示。

```
mysql> DESC INNODB_TRX_ROLLBACK_SEGMENT;

+--------------------+--------------------+------+-----+---------+-------+
| Field              | Type               | Null | Key | Default | Extra |
+--------------------+--------------------+------+-----+---------+-------+
| Segment_id         | bigint(21) unsigned | NO  |     | 0       |       |
| space              | bigint(21) unsigned | NO  |     | 0       |       |
| page_no            | bigint(21) unsigned | NO  |     | 0       |       |
| last_page_no       | bigint(21) unsigned | YES |     | NULL    |       |
| last_offset        | bigint(21) unsigned | NO  |     | 0       |       |
| last_trx_no        | varchar(18)        | NO  |     |         |       |
| update_undo_list   | bigint(21) unsigned | NO  |     | 0       |       |
| update_undo_cached | bigint(21) unsigned | NO  |     | 0       |       |
| insert_undo_list   | bigint(21) unsigned | NO  |     | 0       |       |
| insert_undo_cached | bigint(21) unsigned | NO  |     | 0       |       |
+--------------------+--------------------+------+-----+---------+-------+
10 rows in set (0.00 sec)
```

图 7-16 INNODB_TRX_ROLLBACK_SEGMENT表结构

## 对应事务所使用的回滚段（rollback segment），可以通过

```
mysql>SELECT segment_id,space,page_no
-　FROM INNODB_TRX_ROLLBACK_SEGMENT;
+------------+------+--------+
|segment_id|space|page_no|
+------------+------+--------+
|0|0|6|
|1|0|45|
|2|0|46|
......
128 rows in set(0.00 sec)
```

## 接下来我们来看一下INNODB_TRX_UNDO表，其记录事务对应的undo log，方便DBA和开发人员详细了解每个事务产生的undo量。下面首先根据INNODB_TRX_UNDO表的定义创建测试表t：

```
CREATE TABLE t(
a INT,
b VARCHAR(32),
PRIMARY KEY(a),
KEY(b)
)ENGINE=InnoDB;
```

## 接着插入一条记录，并尝试通过INNODB_TRX_UNDO表来观察该记录对应的undo log的情况：

```
mysql>TBEGIN;
Query OK,0 rows affected(0.00 sec)
mysql>INSERT INTO t SELECT 1,'1';
Query OK,1 row affected(0.00 sec)
```

Records:1 Duplicates:0 Warnings:0

mysql　SELECT*FROM information_schema.INNODB_TRX_UNDO\G;

***************************1.row***************************

trx_id:3001

rseg_id:2

undo_rec_no:0

undo_rec_type:TRX_UNDO_INSERT_REC

size:12

space:0

page_no:334

offset:272

1 row in set(0.00 sec)

---

上述查询结果显示，该事务的ID为3001，rollback segment的ID为2，由于是该事务中的第一个操作，故undo_rec_no为0。由于是插入操作，因此undo类型为TRX_UNDO_INSERT_REC，表示是insert undo log。size表示undo log占用的字节数为12。由于我们的space、page_no、offset表示undo log开始的位置，即在表空间ibdata1中，从页号为334、272偏移量开始的12字节，内容如下所示。

---

01 1c 0b 00 16 04 80 00 00 01 01 10

---

这就是undo log的实际内容。根据之前介绍的undo log格式，可以整理成表所示的内容。

---

01 1c#下一条undo log的位置272+12=0x011c

0b#undo log的类型，TRX_UNDO_INSERT_REC为11

00#undo log的记录编号，undo_rec_no

16#表的ID

04#主键的长度

80 00 00 01#主键的值

01 10#undo log开始的偏移量，272=0x0110

---

下面我们来看看这个undo log。这个rollback segment的ID是2，对应的信息可以从表 INNODB_TRX_ROLLBACK_SEGMENT中获得。这个rollback segment信息 如下所示：

---

mysql＞SELECT segment_id,insert_undo_list,insert_undo_cached

-＞FROM information_schema.INNODB_TRX_ROLLBACK_SEGMENT

-＞WHERE segment_id=2\G;

***************************1.row***************************

segment_id:2

insert_undo_list:1

insert_undo_cached:0

1 row in set(0.00 sec)

---

## 可以看到insert_undo_list为1。接着执行COMMIT动作，然后再观察其变化：

---

mysql＞COMMIT;

Query OK,0 rows affected(0.00 sec)

mysql＞SELECT segment_id,insert_undo_list,insert_undo_cached

-＞FROM information_schema.INNODB_TRX_ROLLBACK_SEGMENT

-＞WHERE segment_id=2\G;

***************************1.row***************************

segment_id:2

insert_undo_list:0

insert_undo_cached:1

1 row in set(0.00 sec)

---

## 此时可以看到insert_undo_list变为0，而insert_undo_cached增加为1。这表示提交的 事务的undo日志放入了链表中，以供下次重用的rollback segment中的undo页链表，可以被 重用。

## 对于上述的delete操作，其对应undo log又是怎么样的呢？

mysql》BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql》DELETE FROM t WHERE a=1;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql》SELECT*FROM information_schema.INNODB_TRX_UNDO\G;

***************************1.row***************************

trx_id:3201

rseg_id:2

undo_rec_no:0

undo_rec_type:TRX_UNDO_DEL_MARK_REC

size:37

space:0

page_no:326

offset:620

1 row in set(0.00 sec)

# 接着通过工具查看，在页面326的偏移量620处可以看到如下记录。

0518260 00 00 00 00 00 00 00 00 00 00 00 00 02 91 0e 00

0518270 16 00 00 00 00 30 01 e0 82 00 00 01 4e 01 10 04

0518280 80 00 00 01 00 0b 00 04 80 00 00 01 03 01 31 02

0518290 6c 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00

## 上述记录的含义是

02 91#下一条undo log的页面偏移量（低位）

0e#undo log类型，TRX_UNDO_DEL_MARK_REC，14

00#undo no

16#table id

00#info bits

00 00 00 30 01 e0#rec□□id

82 00 00 01 4e 01 10#rec□□□□

04#□□□□

80 00 00 01#□□□

00 0b#□□□□□□□

00#□□□□

04#□□□□

80 00 00 01#□□□

03#□□□□□□00□02□□□□

01#□□□□

31#□b□□□□□□□'1'□□□□□

02 6c#□□□□□□□□

# □□rollback segment□□□□□□□□□

mysql□SELECT segment_id,update_undo_list,update_undo_cached

-□FROM information_schema.INNODB_TRX_ROLLBACK_SEGMENT

-□WHERE segment_id=2\G;

***************************1.row***************************

segment_id:2

update_undo_list:1

update_undo_cached:0

1 row in set(0.00 sec)

# □□□□□□□□□□□□□□undo□□□□□cache□□□□□□□□□□□□

mysql□COMMIT;

Query OK,0 rows affected(0.00 sec)

mysql□SELECT segment_id,update_undo_list,update_undo_cached

-□FROM information_schema.INNODB_TRX_ROLLBACK_SEGMENT

-□WHERE segment_id=2\G;

***************************1.row***************************

segment_id:2

update_undo_list:0

update_undo_cached:1

1 row in set(0.00 sec)

---

接下来我们再看一下，当delete一条记录后，删除操作并不会立即进行，而只是设置对应的delete flag的值为1，等待之后的purge操作来进行删除。

因此对update操作，其undo log就是将数据从更新后的1，'1'还原为更新update之前的状态。我们继续通过INNODB_TRX_UNDO观察undo log的信息。

---

mysql□INSERT INTO t SELECT 1,'1';

mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□UPDATE t SET b='2'WHERE a=1;

Query OK,1 row affected(0.00 sec)

Rows matched:1 Changed:1 Warnings:0

mysql□SELECT*FROM information_schema.INNODB_TRX_UNDO\G;

***************************1.row***************************

trx_id:3205

rseg_id:5

undo_rec_no:0

undo_rec_type:TRX_UNDO_UPD_EXIST_REC

size:41

space:0

page_no:318

offset:724

1 row in set(0.00 sec)

---

# 如何解决这个问题呢，对318地址那块，724地址出的数据有这样的内容

04f82d0 00 00 00 00 02 fd 0c 00 16 00 00 00 00 32 04 e0

04f82e0 84 00 00 01 48 01 10 04 80 00 00 01 01 03 01 31

04f82f0 00 0b 00 04 80 00 00 01 03 01 31 02 d4 00 00 00

---

# 这里我们看到

---

02 fd#前一个undo log的结束偏移

0c#undo log类型，TRX_UNDO_UPD_DEL_REC，13

00#undo no

16#table id

00#info bits

00 00 00 32 04 e0#rec trx id

84 00 00 01 48 01 10#rec回滚指针

04#主键长度

80 00 00 01#主键值

01#update vector长度

03#update vector中b的位置

01#update vector中的长度

31#update vector中的值，字符'1'

00 0b#接下来数据的总长度

00#列的位置

04#列的长度

80 00 00 01#列的值

03#列的位置

31#列的值

02 d4#undo log本条记录的开头

# 这样可以很容易地检查某张表在每个事务中产生的回滚段的数量及undo log在系统中的分布情况

mysql□ROLLBACK;

Query OK,1 row affected(0.00 sec)

mysql□UPDATE t SET a=2 WHERE a=1;

Rows matched:1 Changed:1 Warnings:0

mysql□SELECT*FROM information_schema.INNODB_TRX_UNDO

-□ORDER BY undo_rec_no\G;

***************************1.row***************************

trx_id:320F

rseg_id:11

undo_rec_no:0

undo_rec_type:TRX_UNDO_DEL_MARK_REC

size:37

space:0

page_no:324

offset:492

***************************2.row***************************

trx_id:320F

rseg_id:11

undo_rec_no:1

undo_rec_type:TRX_UNDO_INSERT_REC

size:12

space:0

page_no:336

offset:272

2 rows in set(0.00 sec)

□□□□□update□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ TRX_UNDO_DEL_MARK_REC□undo log□□□□□□□□□□□□□□□□□□□□□□□□□□□□ □□□□TRX_UNDO_INSERT_REC□undo log□undo_rec_no□□□□□□□□□□□ □□□□□undo log□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□InnoSQL□□□□□□□□□□□undo□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□ □□□□□□□□undo□□□□

# 7.2.3　purge

delete、update□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□t□□□□□□SQL□□□

```
DELETE FROM t WHERE a=1;
```

□t□□a□□□□□□□□b□□□□□□□□□□□□delete□□□□□□□□□□□□undo log□□□□□□□□□□□□□□□□□□1□□□delete flag□□□□1□□□□□□□□□□□□□□□□□□□□□□B+□□□□□□□□□□□□□□□a□□1□b□□1□□□□□□□□□□□□□□□□□□□□□□undo log□□□□□□□□□□□□□□□□□□□"□□"□□□□□□□purge□□□□□□□□

purge□□□□□□□□delete、update□□□□□□□□□□□□InnoDB□□□□□□□MVCC□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□purge□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□delete□□□□□□purge□□□□□□□□□delete、update□□□□□□□□□□"□□"□□□□□□□□□□□□□□delete□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□undo log□□□□□□□□□□□□□□□□□□□□undo log□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□undo log□□□□□□□□InnoDB□□□□□□□□history□□□□□□□□□□□□□□□□□□undo log□□□□□□□□□□□□□□□□□□

□□7-17□□□□□history list□□□□□□□□□□□□□□□undo log□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□undo page□□□undo log□□□□□□□□□□□□□□□undo page□□□□□□□□□□□□□□□□undo log□trx5□□□□□□□□□□undo log□□□□□□□□□□□□

□□□□purge□□□□□□□InnoDB□□□□□□□□□history list□□□□□□□□□□□□□□□□□□□□□□trx1□□□□□□InnoDB□□□□□□□□□trx1□undo log□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□trx3□□□□□□trx5□□□□□□trx5□□□□□□□□□□□□□□□□□□□□□□□□history list□□□□□□□□□□□□□□□□□□□□trx2□□□□□□trx2□□□□□□□□□□□□□□□□trx6□trx4□□□□□□□□□□□□□undo page2□□□□□□□□□□□□□□□□□□undo page□□□□□□□□□

| History List |
| --- |

trx7 → trx6 → trx5 → trx3 → trx4 → trx2 → trx1

| trx1 | | trx3 | | trx2 |

| trx5 |

| trx6 |

| trx7 |

| trx4 |

undo page1                    undo page2

## 表 7-17 undo log对history链表影响

InnoDB存储引擎的设计中，history list中undo log数量，以及undo page的个数，undo log的数量和对应的回滚，影响着之后介绍的purge的操作。

因为在默认情况下innodb_purge_batch_size的设置，每次purge操作需要清理的undo page个数。在InnoDB1.2版本之前，这个参数的默认值为20，而从1.2版本开始，这个参数的默认值为300。通常来说，该参数设置得越大，每次回收的undo page也就越多，这样可供重用的undo page就越多，减少了磁盘存储空间与分配的开销。不过，该参数设置得太大，则每次需要purge处理更多的undo page，从而导致CPU和磁盘IO过于集中于对undo log的处理，使性能下降。因此对该参数的调整需要由有经验的DBA来操作，并且需要长期观察数据库的运行状态。正如官方的MySQL数据库手册所说的，普通用户不需要调整该参数。

当InnoDB存储引擎的压力非常大时，并不能高效地进行purge操作。那么history list的长度会变得越来越长。全局动态参数innodb_max_purge_lag用来控制history list的长度，若长度大于该参数时，其会"延缓"DML的操作。该参数默认值为0，表示不对history list的长度做任何限制。当大于0时，就会延缓DML的操作，其延缓的算法为：

---

delay=（（length(history_list)-innodb_max_purge_lag）*10)-5

---

delay的单位是毫秒。此外，需要特别注意的是，delay的对象是行，而不是一个DML操作。例如当一个update操作需要更新5行数据时，每行数据的操作都会被delay，故总的延时时间为5*delay。而delay的统计会在每一次purge操作完成后，重新进行计算。

InnoDB1.2版本引入了新的全局动态参数innodb_max_purge_lag_delay，其用来控制delay的最大毫秒数。也就是当上述计算得到的delay值大于该参数时，将delay设置为innodb_max_purge_lag_delay，避免由于purge操作缓慢导致其他SQL线程出现无限制的等待。

# 7.2.4　group commit

若事务为非只读事务，则每次事务提交时需要进行一次fsync操作，以此保证重做日志都已经写入磁盘。当数据库发生宕机时，可以通过重做日志进行恢复。虽然固态硬盘的出现提高了磁盘的性能，然而磁盘的fsync性能是有限的。为了提高磁盘fsync的效率，当前数据库都提供了group commit的功能，即一次fsync可以刷新确保多个事务日志被写入文件。对于InnoDB存储引擎来说，事务提交时会进行两个阶段的操作：

1）修改内存中事务对应的信息，并且将日志写入重做日志缓冲。

2）调用fsync将确保日志都从重做日志缓冲写入磁盘。

步骤2）相对步骤1）是个较慢的过程，这是因为存储引擎需要与磁盘打交道。但当有事务进行这个过程时，其他事务可以进行步骤1）的操作，正在提交的事务完成提交操作后，再次进行步骤2）时，可以将多个事务的重做日志通过一次fsync刷新到磁盘，这样就大大地减少了磁盘的压力，从而提高了数据库的整体性能。对于写入或更新较为频繁的操作，group commit的效果尤为明显。

然而在InnoDB1.2版本之前，在开启二进制日志后，InnoDB存储引擎的group commit功能会失效，从而导致性能的下降。并且在线环境多使用replication环境，因此二进制日志的选项基本都为开启状态，因此这个问题尤为显著。

导致这个问题的原因是在开启二进制日志后，为了保证存储引擎层中的事务和二进制日志的一致性，二者之间使用了两阶段事务，其步骤如下：

1）当事务提交时InnoDB存储引擎进行prepare操作。

2）MySQL数据库上层写入二进制日志。

3）InnoDB存储引擎层将日志写入重做日志文件。

a）修改内存中事务对应的信息，并且将日志写入重做日志缓冲。

b）调用fsync将确保日志都从重做日志缓冲写入磁盘。

一旦步骤2）中的操作完成，就确保了事务的提交，即使在执行步骤3）时数据库发生了宕机。此外，需要注意的是，每个步骤都需要进行一次fsync操作才能保证上下两层数据的一致性。步骤2）的fsync由参数sync_binlog控制，步骤3）的fsync由参数innodb_flush_log_at_trx_commit控制。因此上述整个过程如图7-18所示。

图 7-18　两阶段提交以及InnoDB存储引擎中的各阶段

正如上图，MySQL在开启二阶段提交后，内部会默认使用InnoDB来保证日志文件的一致性，而MySQL为了保证事务提交prepare_commit_mutex，这样会导致以下两个问题。（图中3的问题）：（a）

二阶段提交的过程会被拆解；（b）无法很好的使用group commit功能。

从上图中可以看到，MySQL在开启二阶段提交的时候，会默认使用InnoDB的日志一致性功能，这样会导致备份工具，比如xtrabackup或者ibbackup无法进行正常的热备份replication（图中7-19问题）

图 7-19　InnoDB的两阶段提交。在崩溃恢复时，MySQL对提交做出决定的方式依

赖于存储引擎的排序，但这个协议要求replication。如果在T1之后、提交之前崩溃，则
InnoDB会回滚该事务。如果在T3之后崩溃，则会提交该事务。这个要求的结果就是，
prepare_commit_mutex会序列化提交，强制事务一个接一个地进行group commit，如
同我们在图7-20中所示

图 7-20　移除了prepare_commit_mutex后InnoDB存储引擎可以正常地对MySQL上层进行二进制日志的group commit

那么为什么在2010年，MySQL数据库的一些主要分支，如Facebook MySQL、阿里巴巴的Percona及最早的代码提交者MariaDB的主要开发人员Kristian Nielsen都提出了各自的"补丁"来解决这个问题呢？最后，MySQL数据库官方的版本将实现二进制的group commit，InnoDB存储引擎的group commit也将随之有效，不再需要prepare_commit_mutex锁的存在。二进制日志的组提交，即MySQL 5.6版本开始，为解决这个问题提供了Binary Log Group Commit（BLGC）。

MySQL 5.6 BLGC的实现方式是将事务提交的过程分为几个步骤来完成，如图7-21所示。



图 7-21　MySQL 5.6　BLGC的实现方式

在MySQL数据库上层进行提交时首先按顺序将其放入一个队列中，队列中的第一个事务称为leader，其他事务称为follower，leader控制着follower的行为。BLGC的步骤分为以下三个阶段：

□Flush阶段，将每个事务的二进制日志写入内存中。

□Sync阶段，将内存中的二进制日志刷新到磁盘，若队列中有多个事务，那么仅一次fsync操作就完成了二进制日志的写入，这就是BLGC。

□Commit阶段，leader根据顺序调用存储引擎层事务的提交，由于InnoDB本身就支持group commit，因此解决了原先因为锁prepare_commit_mutex导致group commit失效的问题。

在将事务写入到Commit阶段时，并不用等待其他事务完成Flush阶段，从而group commit一直不停运转。group commit也不是无损的。虽然通过上述的改进大大减少了锁带来的竞争，但并不是没有group commit机制就一定能提高事务的提交性能。

参数binlog_max_flush_queue_time用来控制Flush阶段中等待的时间，即使之前的一组事务完成提交，当前一组的事务也不马上进入Sync阶段，而是至少需要等待一段时间。这样做的好处是group commit的事务数量更多，然而这也可能会导致事务的响应时间变慢。该参数的默认值为0，且推荐设置依然为0，除非用户的MySQL数据库系统中有着大量的连接（如100个连接），并且不断地在进行事务的写入或更新操作。

# 7.3 □□□□□□□

□MySQL□□□□□□□□□□□□□□□□□□□□□□auto commit□□□□□□□SQL□□□□□□□□□□COMMIT□□□□□□□□□□□□□□□□□□□□□□□□BEGIN□START TRANSACTION□□□□□□□□□SET AUTOCOMMIT=0□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□START TRANSACTION|BEGIN□□□□□□□□□□□□□□

□COMMIT□□□□□□□□□□□□□□□□□□□□□□COMMIT□□□□□□□□□□□□□ COMMIT WORK□□□□□□□□□□□□□□COMMIT□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□ROLLBACK□□□□□□□□□□□□□□□□□□□□□□ROLLBACK□□□□□□□□□□□ ROLLBACK WORK□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□SAVEPOINT identifier:SAVEPOINT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SAVEPOINT□

□RELEASE SAVEPOINT identifier□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□ROLLBACK TO[SAVEPOINT]identifier□□□□□□□SAVEPOINT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□UPDATE□□□□□□□□□□□SAVEPOINT□□□□□□□□DELETE□□□□□□□□□DELETE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ROLLBACK TO SAVEPOINT□□□□□□□□□□□□□□□SAVEPOINT□□□□DELETE□□□□□□□□□□UPDATE□□□□□□□□□□□□□

□SET TRANSACTION□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□READ UNCOMMITTED□READ COMMITTED□REPEATABLE READ□SERIALIZABLE□

START TRANSACTION、BEGIN语句都可以显式地开启一个事务。需要注意的是，在MySQL命令行的默认设置下，事务都是自动提交的，BEGIN语句。BEGIN…END语句块只用在存储过程中。而START TRANSACTION后会显式地开启一个事务。

COMMIT、COMMIT WORK语句用来提交事务。但是二者是否完全一致呢？这里的COMMIT WORK语句用来控制事务结束后的CHAIN、RELEASE操作，其中CHAIN的方式和事务链的概念有关。

当参数变量设为completion_type时，这是该参数的默认值，表示0，表示没有任何操作。在这种设置下，COMMIT、COMMIT WORK语句是等价的。当参数completion_type的值为1时，COMMIT WORK等同于COMMIT AND CHAIN，表示立即开启一个相同隔离级别的事务。

---

mysql〉CREATE TABLE t(a INT,PRIMARY KEY(a))ENGINE=INNODB;

Query OK,0 rows affected(0.00 sec)

mysql〉SELECT@@autocommit\G;

***************************1.row***************************

@@autocommit:1

1 row in set(0.00 sec)

mysql〉SET@@completion_type=1;

Query OK,0 rows affected(0.00 sec)

mysql〉BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql〉INSERT INTO t SELECT 1;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql〉COMMIT WORK;

Query OK,0 rows affected(0.01 sec)

mysql〉INSERT INTO t SELECT 2;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql〉INSERT INTO t SELECT 2;

ERROR 1062(23000):Duplicate entry'2'for key'PRIMARY'

mysql□ROLLBACK;

Query OK,0 rows affected(0.00 sec)

#□□□□□□□□□1□□□□□□□□□2□□□□

mysql□SELECT*FROM t\G;

***************************1.row***************************

a:1

1 row in set(0.00 sec)

---

□□□□□□□□□□□□completion_type□1□□□□□□□□COMMIT WORK□□□□1□□□□□□□□□□□□□□2□□□□□□□□BEGIN□□□□START TRANSACTION□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□2□□□□□□□□□□□□□ROLLBACK□□□□□□□□□□□1□□□□□□□2□□□□□□□□□□completion_type□1□□COMMIT WORK□□□□□□□□□□□□□□□□□INSERT INTO t SELECT 2□□□□□□□□□□□□□□□□□□□□□□2□□□□□□□□□□□□□t□□□

□□completion_type□2□□□COMMIT WORK□□□□COMMIT AND RELEASE□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□SET@@completion_type=2;

Query OK,0 rows affected(0.00 sec)

mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO t SELECT 3;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql□COMMIT WORK;

Query OK,0 rows affected(0.01 sec)

mysql□SELECT@@version\G;

ERROR 2006(HY000):MySQL server has gone away

No connection.Trying to reconnect...

Connection id:54

Current database:test

***************************1.row***************************

@@version:5.1.45-log

1 row in set(0.00 sec)

---

□□□□□□□□□□□□□□□□completion_type□□□2□□□COMMIT WORK□□□□□□□□□□SELECT@@version□□□□ERROR 2006□HY000□□MySQL server has gone away□□□□□□□□□□□□□□□□□□□□□□□□□□□□COMMIT WORK□□□□□□□□□□□□□□□□

ROLLBACK□ROLLBACK WORK□COMMIT□COMMIT WORK□□□□□□□□□□□□□□□□□

SAVEPOINT□□□□□□□□□□□□□□□□ROLLBACK TO SAVEPOINT□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□ROLLBACK TO SAVEPOINT t1;

ERROR 1305(42000):SAVEPOINT t1 does not exist

---

InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□CREATE TABLE t(a INT,PRIMARY KEY(a))ENGINE=INNODB;

Query OK,0 rows affected(0.00 sec)

mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO t SELECT 1;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql□INSERT INTO t SELECT 1;

ERROR 1062(23000):Duplicate entry'1'for key'PRIMARY'

mysql□SELECT*FROM t\G;

***************************1.row***************************

a:1

1 row in set(0.00 sec)

---

可以看到，插入的值为1的记录，因为主键重复抛出了1062错误。但是可以发现，事务其实并没有进行回滚，如果想让整个事务进行回滚，必须执行COMMIT、ROLLBACK等命令。

这同样也符合上面所说的ROLLBACK TO SAVEPOINT，虽然有ROLLBACK的字样，但其并不是真正地进行事务的回滚，而是在保存点ROLLBACK TO SAVEPOINT，事务并没有结束。想要完整地回滚事务，需要COMMIT、ROLLBACK等命令。

---

mysql□CREATE TABLE t(a INT,PRIMARY KEY(a))ENGINE=INNODB;

Query OK,0 rows affected(0.00 sec)

mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO t SELECT 1;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql□SAVEPOINT t1;

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO t SELECT 2;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql□SAVEPOINT t2;

Query OK,0 rows affected(0.00 sec)

mysql□RELEASE SAVEPOINT t1;

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO t SELECT 2;

ERROR 1062(23000):Duplicate entry'2'for key'PRIMARY'

mysql□ROLLBACK TO SAVEPOINT t2;

Query OK,0 rows affected(0.00 sec)

mysql□SELECT*FROM t;

+---+

|a|

+---+

|1|

|2|

+---+

2 rows in set(0.00 sec)

mysql□ROLLBACK;

Query OK,0 rows affected(0.00 sec)

mysql□SELECT*FROM t;

Empty set(0.00 sec)

---

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□ROLLBACK TO SAVEPOINT t2□□□□□□□□□□t2□□□□□□□□□□□□□□□□□□□□ROLLBACK□□□□□□□□□□□□□□□□□□□□□□□ROLLBACK TO SAVEPOINT□□□□□□□□□□□□□

# 7.4　隐式提交的SQL语句

以下这些SQL语句会产生一个隐式的提交操作，即执行完这些语句后，会有一个隐式的COMMIT操作。

□DDL语句：ALTER DATABASE...UPGRADE DATA DIRECTORY NAME、ALTER EVENT、ALTER PROCEDURE、ALTER TABLE、ALTER VIEW、CREATE DATABASE、CREATE EVENT、CREATE INDEX、CREATE PROCEDURE、CREATE TABLE、CREATE TRIGGER、CREATE VIEW、DROP DATABASE、DROP EVENT、DROP INDEX、DROP PROCEDURE、DROP TABLE、DROP TRIGGER、DROP VIEW、RENAME TABLE、TRUNCATE TABLE。

□用来隐式地修改MySQL架构的操作：CREATE USER、DROP USER、GRANT、RENAME USER、REVOKE、SET PASSWORD。

□管理语句：ANALYZE TABLE、CACHE INDEX、CHECK TABLE、LOAD INDEX INTO CACHE、OPTIMIZE TABLE、REPAIR TABLE。

注意　对于Microsoft SQL Server的用户来说，通常不会有隐式提交的操作，DDL也是可以回滚的。但在Microsoft SQL Server中，即使支持DDL的事务，个人也不建议将DDL和DML放在同一个事务中，这是因为DDL的回滚是比较昂贵的操作。

还有一点需要注意的是，TRUNCATE TABLE语句是DDL，因此虽然和对整张表执行DELETE的结果是一样的，但它是不能被回滚的（这又和Microsoft SQL Server数据库有所不同）。

---

```
mysql〉SELECT*FRM t\G;

**************************1.row**************************

a:1

**************************2.row**************************

a:2

2 rows in set(0.00 sec)

mysql〉BEGIN;

Query OK,0 rows affected(0.01 sec)
```

mysql›TRUNCATE TABLE t;

Query OK,0 rows affected(0.00 sec)

mysql›ROLLBACK;

Query OK,0 rows affected(0.00 sec)

mysql›SELECT*FROM t;

Empty set(0.00 sec)

# 7.5 □□□□□□□□□□□

□□InnoDB□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□
□Question Per Second□QPS□□□□□□□□□□□□□□□□□□□□□□□
□Transaction Per Second□TPS□□

□□TPS□□□□□□com_commit+com_rollback□/time□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□autocommit=1□□□
□□□□□□com_commit□com_rollback□□□□□□□□

---

```
mysql□SHOW GLOBAL STATUS LIKE'com_commit'\G;

***************************1.row***************************

Variable_name:Com_commit

Value:5

1 row in set(0.00 sec)

mysql□INSERT INTO t SELECT 3;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql□SELECT*FROM t\G;

***************************1.row***************************

a:1

***************************2.row***************************

a:2

***************************3.row***************************

a:3

3 rows in set(0.00 sec)

mysql□SHOW GLOBAL STATUS LIKE'com_commit'\G;

***************************1.row***************************

Variable_name:Com_commit

Value:5
```

MySQL数据库中另外两个监控变量handler_commit、handler_rollback对于事务的监控意义不大，因为在MySQL 5.1版本中，这两个值由存储引擎InnoDB插件提供，并非都由InnoDB存储引擎本身提供，故这两个参数在MySQL数据库的"正常"运行中可能不发生任何变化。而对于目前大多数的OLTP应用，上述的com_commit和com_rollback却又关系到每秒请求数，故对这两个参数的监控还是非常有意义的。

# 7.6　事务的隔离级别

事务可以设置不同的隔离级别，隔离级别越低，并发程度越高，性能越好，但是同时也可能会出现更多的数据一致性问题。在ISO和ANIS SQL标准中制定了四种事务的隔离级别，但是不同的数据库厂商也可能支持不同的隔离级别，例如Oracle数据库就不支持READ UNCOMMITTED和REPEATABLE READ这两种隔离级别。

SQL标准定义的四种隔离级别为：

❑READ UNCOMMITTED

❑READ COMMITTED

❑REPEATABLE READ

❑SERIALIZABLE

READ UNCOMMITTED被称为浏览访问（browse access），仅仅针对事务而言的。READ COMMITTED被称为游标稳定（cursor stability）。REPEATABLE READ是2.9999°的隔离，没有幻读的保护。SERIALIZABLE被称为隔离，或3°的隔离。SQL和SQL2标准的默认事务隔离级别是SERIALIZABLE。

InnoDB存储引擎默认支持的隔离级别是REPEATABLE READ，但是与标准SQL不同的是，InnoDB存储引擎在REPEATABLE READ事务隔离级别下，使用Next-Key Lock锁的算法，因此避免幻读的产生。这与其他数据库系统（如Microsoft SQL Server数据库）是不同的。所以说，InnoDB存储引擎在默认的REPEATABLE READ的事务隔离级别下已经能完全保证事务的隔离性要求，即达到SQL标准的SERIALIZABLE隔离级别。

隔离级别越低，事务请求的锁越少或保持锁的时间就越短。这也是为什么大多数数据库系统默认的事务隔离级别是READ COMMITTED。

虽然数据库厂商设置的隔离级别是SERIALIZABLE，但是我们在平时工作中，更多出于对性能的考虑（Jim Gray在《Transaction Processing》一书中提到，每提高一个隔离级别，性能可能下降50%），设置使用的隔离级别是SERIALIZABLE的情况并不多见!!!而是对InnoDB存储引擎默认的REPEATABLE READ级别有一种"不安全感"，大多数数据库系统都使用READ COMMITTED作为事务默认的隔离级别，一般交易也够用了。

为InnoDB设置隔离级别，可以使用下面的语句。其中的级别可以是下面的一种。

```
SET[GLOBAL|SESSION]TRANSACTION ISOLATION LEVEL

{

READ UNCOMMITTED

|READ COMMITTED

|REPEATABLE READ

|SERIALIZABLE

}
```

也可以在MySQL的配置文件中设置隔离级别。具体做法是，修改MySQL的配置文件，在[mysqld]节中加入下面的行。

```
[mysqld]

transaction-isolation=READ-COMMITTED
```

在当前会话中查询隔离级别的方法如下。

```
mysql>SELECT@@tx_isolation\G;

***************************1.row***************************

@@tx_isolation:REPEATABLE-READ

1 row in set(0.01 sec)
```

在全局范围内查询隔离级别的方法如下。

```
mysql>SELECT@@global.tx_isolation\G;

***************************1.row***************************

@@global.tx_isolation:REPEATABLE-READ

1 row in set(0.00 sec)
```

在SERIALIABLE的隔离级别下，InnoDB会对每个读取操作对SELECT语句自动加上
LOCK IN SHARE MODE，即为每个读取操作加一个共享锁。因此在这个隔离级别下，读占用了锁，对一致性的非锁定读不再予以支持。SERIALIZABLE的事务隔离级别主要用于分布式事务，即符合well-formed的两段锁two-phrased锁提交的分布式事务处理。

由于InnoDB默认采用REPEATABLE READ的隔离级别，而符合3°的隔离要求，因此一般不在本地事务中采用SERIALIABLE的隔离级别。SERIALIABLE的隔离级别带来的开销会降低InnoDB数据库的并发处理能力。

在READ COMMITTED的事务隔离级别下，除了唯一性的约束检查以及外键约束的检查需要gap lock，InnoDB不会使用gap lock的锁算法。然而这个设置会对于数据库的恢复与复制带来一定的影响。在MySQL 5.1时，READ COMMITTED的隔离级别只能工作在replication二进制日志格式为ROW的情况下。如果二进制日志格式为STATEMENT，则会抛出如下的一个异常。

---

mysql□CREATE TABLE a(

-□b INT,PRIMARY KEY(b)

-□)ENGINE=INNODB;

Query OK,0 rows affected(0.01 sec)

mysql□SET@@tx_isolation='READ-COMMITTED';

Query OK,0 rows affected(0.00 sec)

mysql□SELECT@@tx_isolation\G;

***************************1.row***************************

@@tx_isolation:REPEATABLE-READ

1 row in set(0.00 sec)

mysql□BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql□INSERT INTO a SELECT 1;

ERROR 1598(HY000):Binary logging not possible.Message:Transaction level'READ-COMMITTED'in InnoDB is not safe for binlog mode'STATEMENT'

---

# 在MySQL 5.0及更高版本中，在ROW模式下，可将事务隔离级别设置为READ COMMITTED，而不是innodb_locks_unsafe_for_binlog设为1，它会有如下的STATEMENT、设置READ COMMITTED都会有所限制：

mysql＞SELCT@@version\G

***************************1.row***************************

@@version:5.0.77-log

1 row in set(0.00 sec)

mysql＞SHOW VARIABLES LIKE'innodb_locks_unsafe_for_binlog'\G;

***************************1.row***************************

Variable_name:innodb_locks_unsafe_for_binlog

Value:ON

1 row in set(0.00 sec)

mysql＞SET@@tx_isolation='READ-COMMITTED';

Query OK,0 rows affected(0.00 sec)

mysql＞BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql＞INSERT INTO a SELECT 1;

Query OK,0 rows affected(0.00 sec)

mysql＞COMMIT;

Query OK,0 rows affected(0.00 sec)

mysql＞SELECT*FROM a\G;

***************************1.row***************************

b:1

***************************2.row***************************

b:2

***************************3.row***************************

b:4

***************************4.row***************************

b:5

4 rows in set(0.00 sec)

# 步骤二：master上开启一个事务A，删除表中键值大于等于某数的

#Session A on master

mysql，BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql，DELETE FROM a WHERE b，=5;

Query OK,4 rows affected(0.01 sec)

# 步骤三：master上再开启一个事务B，向表中插入键值等于某数的记

#Session B on master

mysql，BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql，INSERT INTO a SELECT 3;

Query OK,0 rows affected(0.01 sec)

mysql，COMMIT;

Query OK,0 rows affected(0.00 sec)

# 步骤四：A事务提交，并查询a表的结果。

#Session A on master

mysql，COMMIT;

Query OK,0 rows affected(0.00 sec)

mysql，SELECT*FROM a\G;

***************************1.row***************************

b:3

可以看到slave上没有插入记录。

#Slave

mysql␣SELECT*FROM a;

Empty set(0.00 sec)

___

那么为什么会发生这样奇怪的现象呢?这里的问题出在:

①在READ COMMITTED事务隔离级别下,事务没有使用gap lock进行加锁,因此用户在会话B中可以在事务完成5的插入操作之前插入了一条记录;

②STATEMENT格式记录的是master上产生的SQL语句,因此在master服务器上用户会话A的事务先完成,由STATEMENT格式记录,插入的数据会被传到到子服务器上进行重放。

这个问题的解决方法是使用同步的二进制日志格式,所以强烈建议使用READ REPEATABLE的事务隔离级别。其实前面已经论述过,在这种情况下,master和slave数据会发生不一致。

在MySQL 5.1版本之后,由于支持了ROW格式的二进制日志记录格式,上述问题不再存在,用户可以放心地在READ COMMITTED隔离级别下使用。但是在READ COMMITTED隔离级别下,用户使用二进制日志时应该选择ROW格式。而且这个设置应该是默认的设置,否则复制变得没有意义。这也是默认推荐事务隔离级别设置的原因。InnoDB存储引擎创始人HeikkiTuuri在http:‖bugs.mysql.com/bug.php?id=33210案例中也强调了使用ROW格式二进制日志的重要性。

# 7.7 分布式事务

## 7.7.1 MySQL数据库分布式事务

InnoDB存储引擎提供了对XA事务的支持，并通过XA事务来支持分布式事务的实现。分布式事务指的是允许多个独立的事务资源（transactional resources）参与到一个全局的事务中。事务资源通常是关系型数据库系统，但也可以是其他类型的资源。全局事务要求在其中的所有参与的事务要么都提交，要么都回滚，这对于事务原有的ACID要求又有了提高。另外，在使用分布式事务时，InnoDB存储引擎的事务隔离级别必须设置为SERIALIZABLE。

XA事务允许不同数据库之间的分布式事务，如一台服务器是MySQL数据库的，另一台是Oracle数据库的，又可能还有一台服务器是SQL Server数据库的，只要参与在全局事务中的每个节点都支持XA事务。分布式事务可能在银行系统的转账中比较常见，如用户David需要从上海的账户转10 000元到北京的用户Mariah的账户：

---

```
#Bank@Shanghai：

UPDATE account SET money=money-10000 WHERE user='David';

#Bank@Beijing

UPDATE account SET money=money+10000 WHERE user='Mariah';
```

---

在这种情况下，一定需要使用分布式事务来保证数据的安全。如果用户可以接受两个误差内的数据最终一致性，如上海的用户David是减去了10 000元，但是北京的用户Mariah没有立即收到这10 000元，在有且仅有一个误差的情况下两个数据是一致的。可以最后发现上海的用户David减去了10 000元，但是北京的用户Mariah没有收到这10 000元。

XA事务由一个或多个资源管理器（Resource Managers）、一个事务管理器（Transaction Manager）以及一个应用程序（Application Program）组成。

□资源管理器：提供访问事务资源的方法。通常一个数据库就是一个资源管理器。

□事务管理器：协调参与全局事务中的各个事务。需要和参与全局事务的所有资源管理器进行通信。

□应用程序：定义事务的边界，指定全局事务中的操作。

在MySQL数据库中的分布式事务，应用程序由一个MySQL客户端或者多个、位于多个MySQL服务器上的客户端组成，如图7-22所示为一个典型的分布式事务的模型。



图 7-22   分布式事务模型

分布式事务采用两段式提交（two-phase commit）的方式。第一阶段所有的事务节点开始准备，告诉事务管理器它们准备好提交了。第二阶段事务管理器告诉每个节点是否要提交。在这个时候可能发生某个节点执行ROLLBACK，或者COMMIT。在分布式事务环境下，应用程序可以同时操作多个不同的数据库节点。

分支事务必须等到接收到全局事务的PREPARE消息时才能进入准备提交的阶段，再根据全局事务发送的COMMIT或者ROLLBACK来决定。

## MySQL数据库XA语句的SQL语法如下：

---

XA{START|BEGIN}xid[JOIN|RESUME]

XA END xid[SUSPEND[FOR MIGRATE]]

XA PREPARE xid

XA COMMIT xid[ONE PHASE]

XA ROLLBACK xid

XA RECOVER

---

# 单个节点上运行XA事务的过程：

---

mysql》XA START'a';

Query OK,0 rows affected(0.00 sec)

mysql》INSERT INTO z SELECT 11;

Query OK,1 row affected(0.00 sec)

Records:1 Duplicates:0 Warnings:0

mysql》XA END'a';

Query OK,0 rows affected(0.00 sec)

mysql》XA PREPARE'a';

Query OK,0 rows affected(0.05 sec)

mysql》XA RECOVER\G;

***************************1.row***************************

formatID:1

gtrid_length:1

bqual_length:0

data:a

1 row in set(0.00 sec)

mysql口XA COMMIT'a';

Query OK,0 rows affected(0.05 sec)

---

口口口口口口口口口口口口口口口口口口口口口口口口MySQL口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口Java口JTA口Java Transaction API口口口口口口MySQL口口口口口口口口口口口口口口口口口口口口口口口口口口口API口口口口口口口口口口口口口口口JTA口口口MySQL口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口口

---

```java
import java.sql.Connection;

import javax.sql.XAConnection;

import javax.transaction.xa.*;

import com.mysql.jdbc.jdbc2.optional.MysqlXADataSource;

import java.sql.*;

class MyXid implements Xid

{

public int formatId;

public byte gtrid[];

public byte bqual[];

public MyXid(){

}

public MyXid(int formatId,byte gtrid[],byte bqual[])

{

this.formatId=formatId;

this.gtrid=gtrid;

this.bqual=bqual;

}

public int getFormatId()

{

return formatId;

}
```

```java
public byte[]getBranchQualifier()

{

return bqual;

}

public byte[]getGlobalTransactionId()

{

return gtrid;

}

}

public class xa_demo{

public static MysqlXADataSource GetDataSource(

String connString,

String user,

String passwd){

try{

MysqlXADataSource ds=new MysqlXADataSource();

ds.setUrl(connString);

ds.setUser(user);

ds.setPassword(passwd);

return ds;

}

catch(Exception e){

System.out.println(e.toString());

return null;

}

}

public static void main(String[]args){

String connString1="jdbc:mysql://192.168.24.43:3306/bank_shanghai";

String connString2="jdbc:mysql://192.168.24.166:3306/bank_

beijing";
```

```java
try{

MysqlXADataSource ds1=

GetDataSource(connString1,"peter","12345");

MysqlXADataSource ds2=

GetDataSource(connString2,"david","12345");

XAConnection xaConn1=ds1.getXAConnection();

XAResource xaRes1=xaConn1.getXAResource();

Connection conn1=xaConn1.getConnection();

Statement stmt1=conn1.createStatement();

XAConnection xaConn2=ds2.getXAConnection();

XAResource xaRes2=xaConn2.getXAResource();

Connection conn2=xaConn2.getConnection();

Statement stmt2=conn2.createStatement();

Xid xid1=new MyXid(

100,

new byte[]{0x01},

new byte[]{0x02});

Xid xid2=new MyXid(

100,

new byte[]{0x11},

new byte[]{0x12});

try{

xaRes1.start(xid1,XAResource.TMNOFLAGS);

stmt1.execute("

UPDATE account SET money=money-10000

WHERE user='david'"

);

xaRes1.end(xid1,XAResource.TMSUCCESS);

xaRes2.start(xid2,XAResource.TMNOFLAGS);

stmt2.execute("
```

```
UPDATE account SET money=money+10000

WHERE user='mariah'"

);

xaRes2.end(xid2,XAResource.TMSUCCESS);

int ret2=xaRes2.prepare(xid2);

int ret1=xaRes1.prepare(xid1);

if(ret1==XAResource.XA_OK

　　ret2==XAResource.XA_OK){

xaRes1.commit(xid1,false);

xaRes2.commit(xid2,false);

}

}catch(Exception e){

e.printStackTrace();

}

}catch(Exception e){

System.out.println(e.toString());

}

}

}
```

# 说明：innodb_support_xa用来开启或关闭分布式XA事务的支持。默认值为ON表示

mysql　SHOW VARIABLES LIKE'innodb_support_xa'\G;

***************************1.row***************************

Variable_name:innodb_support_xa

Value:ON

1 row in set(0.01 sec)

# 7.7.2 内部XA事务

内部事务主要指的是在单个实例进行的分布式事务，MySQL中最典型的就是MySQL内部事务，即在单个实例上协调存储引擎和二进制日志之间的内部事务，当然还包括内部事务建立之上的XA事务。

在单个实例上使用XA事务协调binlog和InnoDB存储引擎事务的提交，即保证事务提交前写入的二进制日志binlog被刷新到磁盘，另外保证在存储引擎InnoDB层提交的事务同样被记录在二进制日志中，因为二进制日志记录了数据库的逻辑变化，当InnoDB层已经提交的事务，当从机slave同步时依赖master的二进制日志来进行，这就要求事务提交前将其写入二进制日志，如图7-23所示。



图 7-23 主从结构replication过程中的事务处理

在图7-23中表示，先执行①，②步骤，再③，因MySQL主从同步架构要求将事务提交前写入二进制日志，因此，MySQL先进行binlog，InnoDB事务的协调采用XA事务来进行，先在InnoDB存储引擎层进行PREPARE操作，并写入xid信息，记录此时的日志逻辑序号。

如7-24所示，当底层InnoDB存储引擎需要和MySQL数据库层的二进制日志进行同步，因为MySQL数据库上层还有一个协调者，UXID的值一般可认为是二进制日志中的偏移量信息。分布式事务



图　7-24　MySQL数据库层面上的XA事务提交时分布式事务处理

# 7.8 长事务

## 7.8.1 长事务的问题

回顾一下本章开头介绍的一个修改用户表中记录金额的例子，例子中会将用户分成不

---

```
CREATE PROCEDURE load1(count INT UNSIGNED)

BEGIN

DECLARE s INT UNSIGNED DEFAULT 1;

DECLARE c CHAR(80)DEFAULT REPEAT('a',80);

WHILE s□=count DO

INSERT INTO t1 SELECT NULL,c;

COMMIT;

SET s=s+1;

END WHILE;

END;
```

---

在上述的存储过程中，每插入一条记录就进行一次COMMIT操作。而我们知道，InnoDB存储引擎默认为自动提交，所以这里显式地加上了COMMIT操作。这样做的目的是为了大量插入时不至于爆发很大的回滚日志。

---

```
CREATE PROCEDURE load2(count INT UNSIGNED)

BEGIN

DECLARE s INT UNSIGNED DEFAULT 1;

DECLARE c CHAR(80)DEFAULT REPEAT('a',80);

WHILE s□=count DO

INSERT INTO t1 SELECT NULL,c;

SET s=s+1;

END WHILE;
```

END;

---

下面创建第三个存储过程，它和第二个存储过程几乎完全一样，但它使用了事务处理。它插入了10 000行记录，但每插入5000行记录就提交一次。我们每隔5000行记录就提交一次，原因是希望减少事务处理的规模。下面是load3存储过程的代码，它的关键部分是INSERT语句前后的几行代码：

---

```
CREATE PROCEDURE load3(count INT UNSIGNED)

BEGIN

DECLARE s INT UNSIGNED DEFAULT 1;

DECLARE c CHAR(80)DEFAULT REPEAT('a',80);

START TRANSACTION;

WHILE s<=count DO

INSERT INTO t1 SELECT NULL,c;

SET s=s+1;

END WHILE;

COMMIT;

END;
```

---

# 对这3个存储过程进行测试，如下：

---

```
mysql>CALL load1(10000);

Query OK,0 rows affected(1 min 3.15 sec)

mysql>TRUNCATE TABLE t1;

Query OK,0 rows affected(0.05 sec)

mysql>CALL load2(10000);

Query OK,1 row affected(1 min 1.69 sec)

mysql>TRUNCATE TABLE t1;

Query OK,0 rows affected(0.05 sec)

mysql>CALL load3(10000);
```

Query OK,0 rows affected(0.63 sec)

---

接着我们来观察最为有趣的一个例子!在看下面的例子前，首先清空表中数据，调用load1、load2各插入10 000条数据，然后观察在调用load3更新数据时，只输入1万分之一的更新数据（调用load2），看看会给数据更新性能带来什么影响（load3更新数据）。

---

mysql：BEGIN;

Query OK,0 rows affected(0.00 sec)

mysql：CALL load2(10000);

Query OK,1 row affected(0.56 sec)

mysql：COMMIT;

Query OK,0 rows affected(0.03 sec)

---

读者可能对这个实验中遇到的问题感到疑惑，因为在InnoDB存储引擎中，事务的隔离级别是可重复读，在这个隔离级别下，读取的数据不会因为其他事务的更新而发生变化。这就类似于在Oracle数据库中，如果回滚段undo中的Snapshot Too Old。但是在MySQL的InnoDB存储引擎中，由于采用了多版本并发控制的机制，所以不会出现这样的问题，这也是其性能优异的一个重要原因。

# 7.8.2 　□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□autocommit□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□SET autocommit=0;

Query OK,0 rows affected(0.00 sec)

---

□□□□□□START TRANSACTION□BEGIN□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□completion_type□□0□□MySQL□□□□□□□□SET AUTOCOMMIT=0□□□□□□□COMMIT□ROLLBACK□□□□□□□□□□□□SET AUTOCOMMIT=1□

□□□□□□□□□□□□□□API□□□□□□□□□□□□□MySQL C API□□□□□□□□□□□□□□□□□□□□MySQL Python API□□□□□□□SET AUTOCOMMIT=0□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□API□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 7.8.3 □□□□□□

InnoDB□□□□□□□□□□□□□HANDLER□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

```
CREATE PROCEDURE sp_auto_rollback_demo()

BEGIN

DECLARE EXIT HANDLER FOR SQLEXCEPTION ROLLBACK;

START TRANSACTION;

INSERT INTO b SELECT 1;

INSERT INTO b SELECT 2;

INSERT INTO b SELECT 1;

INSERT INTO b SELECT 3;

COMMIT;

END;
```

□□□□□sp_auto_rollback_demo□□□□□□□□□exit□□□□HANDLER□□□□□□□□□□□□□□□□□□□□□

```
mysql□SHOW CREATE TABLE b\G;

***************************1.row***************************

Table:b

Create Table:CREATE TABLE'b'(

'a'int(11)NOT NULL DEFAULT'0',

PRIMARY KEY('a')

)ENGINE=InnoDB DEFAULT CHARSET=latin1

1 row in set(0.00 sec)
```

## 存储过程成功返回值1，并且由于数据中含有重复值，所以数据没有被插入到表中，如下所示：

mysql>CALL sp_auto_rollback_demo;

Query OK,0 rows affected(0.06 sec)

mysql>SELECT*FROM b;

Empty set(0.00 sec)

## 如果错误发生在事务提交之后，存储过程sp_auto_rollback_demo将正常返回。如果错误发生在事务提交之前，存储过程sp_auto_rollback_demo将回滚事务，返回错误码。下面是一个简单的示例代码如下所示：

```
CREATE PROCEDURE sp_auto_rollback_demo()

BEGIN

DECLARE EXIT HANDLER FOR SQLEXCEPTION BEGIN ROLLBACK;SELECT-1;END;

START TRANSACTION;

INSERT INTO b SELECT 1;

INSERT INTO b SELECT 2;

INSERT INTO b SELECT 1;

INSERT INTO b SELECT 3;

COMMIT;

SELECT 1;

END;
```

## 从下面的代码可以看到，返回值-1，表明存储过程执行失败。如果返回值1，就表明存储过程执行成功，如下所示：

mysql>CALL sp_auto_rollback_demo()\G;

***************************1.row***************************

-1:-1

1 row in set(0.04 sec)

mysql>SELECT*FROM b;

```
Empty set(0.00 sec)
```

---

一个存储过程中包含多个数据库操作语句时，如果其中的某个数据库操作语句执行失败了，我们希望整个存储过程都回滚。

在这一点上，我们必须强调一下，在Microsoft SQL Server中，我们可以通过在Microsoft SQL Server中设置参数（SET XABORT ON）来实现如果某条语句在Microsoft SQL Server中执行失败了，则整个事务都自动回滚。遗憾的是，在Microsoft SQL Server，在MySQL中并没有这样的参数可以设置。

我们必须通过在存储过程中使用BEGIN、COMMIT、ROLLBACK来实现整个存储过程的回滚。下面让我们来看一个具体的例子。我们使用Python来创建两个存储过程，它们分别是sp_rollback_demo，与其不同的是，sp_rollback_demo的第二版本，sp_auto_rollback_demo。这个两存储过程的定义如下所示：

---

```
CREATE PROCEDURE sp_rollback_demo()

BEGIN

INSERT INTO b SELECT 1;

INSERT INTO b SELECT 2;

INSERT INTO b SELECT 1;

INSERT INTO b SELECT 3;

END;
```

---

而sp_auto_rollback_demo存储过程的定义与sp_rollback_demo的定义则完全不同，在它的定义中我们使用了如下所示的方式。我们将代码保存为test_demo.py，如下所示：

---

```
#!/usr/bin/env python

#encoding=utf-8

import MySQLdb

try:

conn=

MySQLdb.connect(host="192.168.8.7",user="root",passwd="xx",db="test")
```

```
cur=conn.cursor()

cur.execute("SET autocommit=0")

cur.execute("CALL sp_rollback_demo")

cur.execute("COMMIT")

except Exception,e:

cur.execute("ROLLBACK")

print e
```

## 接下来执行test_demo.py，并观察输出结果：

```
[root@nineyou0-43 ]#python test_demo.py

starting rollback

(1062,"Duplicate entry'1'for key'PRIMARY'")
```

从结果上可以看到，存储过程在被应用程序调用的时候发生了错误，返回了一个错误代码1062，错误代码所对应的错误信息是Duplicate entry'1'for key'PRIMARY'，说明操作发生了主键冲突，同时应用程序对操作进行了回滚。

# 7.9 长事务

长事务(Long-Lived Transactions)指的是那些运行相当长一段时间的事务，这类事务违背了事务简短的设计理念，却又在生产环境中客观存在。例如，银行系统中的计算账户1年利息，对account表进行更新的操作：

---

```
UPDATE account

SET account_total=account_total+(1+interest_rate)
```

---

如果这样的语句要对全国上千万个账户计算1年利息，可能需要运行4～5个小时，这将是一个对数据库DBA来说十分痛苦的事务。这类长事务带来的问题是，由于ACID特性，数据库可能需要持有大量的锁资源，同时回滚段也会变得巨大无比。一旦运行这类事务的过程中数据库发生了问题，数据库需要花费和事务执行相同甚至更长的时间进行恢复操作，这是难以接受的。因此，对于长事务的问题，有时可以通过转化为小批量(mini batch)的事务来进行处理。当事务发生错误时，只需要回滚一部分数据，然后接着上次已完成的事务继续进行。

例如可以通过如下的代码来完成上述的计算利息的任务，从而将一个需要处理上千万个账户的大事务，分解为每次处理一万个账户的小事务：

---

```
void ComputeInterest（double interest_rate){

long last_account_done,max_account_no,log_size;

int batch_size=100000;

EXEC SQL SELECT COUNT(*)INTO log_size FROM batchcontext;

if(SQLCODE!=0||log_size==0){

EXEC SQL DROP TABLE IF EXISTS batchcontext;

EXEC SQL CREATE TABLE batchcontext(last_account_done BIGINT);

last_account_done=0;

INSERT INTO batchcontext SELECT 0;

}

else{

EXEC SQL SELECT last_account_no
```

```
INTO last_account_done

FROM batchcontext;

}

EXEC SQL SELECT COUNT(*)INTO max_account_no

FROM account LOCK IN SHARE MODE;

WHILE(last_account_no□max_account_no){

EXEC SQL START TRANSACTION;

EXEC SQL UPDATE account

SET account_total=account_total*(1+interest_rate);

WHERE account_no

BETWEEN last_account_no

AND last_account_no+batch_size;

EXEC SQL UPDATE batchcontext

SET last_account_done=last_account_done+batch_size;

EXEC SQL COMMIT WORK;

last_account_done=last_account_done+batch_size;

}

}
```

---

这段代码中的编号为1的事务只需要不超过最初10个账户的锁就足够了。这里需要进行一点说明，事务改变了批处理上下文表batchcontext的信息，存储了最后一次被处理的账户ID，当系统发生故障的时候，我们可以从它存储的账户ID所指示的那个断点位置进行恢复操作，而不必回退到最开始的位置。通过对批处理上下文表batchcontext的观察我们也可以容易地了解批处理操作的进展情况，例如，假设编号为1的事务运行之后，batchcontext表存储的账户ID为4000，那么我们就可以知道大概有40%的账户

已经被处理了，因为在账户表account中，max_account_no的值为最大的账户编号，它给出了所有账户的数目，也就表明了所有账户的总体进度处理情况。

# 7.10　□□

□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□ACID□□□□□Atomicity(□□□)□Consistency(□□□)□Isolation(□□□)□Durability(□□□)□□□□□□□□□6□□□□□□□□□□□□□□□□□□□□□□□□□redo□undo□□□□□□□□redo□undo□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□REPEATABLE READ□□□□□□SQL□□□□□□□□□□□□□□□□InnoDB□□□□□REPEATABLE READ□□□□□□□□□□□3°□□□□□□□

□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□——□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□START TRANSACTION□COMMIT□ROLLBACK□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL InnoDB□□□□□□□□□

# 第8章 备份与恢复

　　作为DBA，最担心的莫过于数据丢失，因为数据一旦丢失，可能会对企业造成难以估量的影响。尽管RAID技术可以保证在磁盘损坏的情况下数据不丢失，但它并不能保证数据逻辑上的正确性，因此DBA平时最重要的工作内容之一就是对数据进行备份。

　　在第2章中介绍了InnoDB的工作方式，在了解了整个MySQL数据库系统后，本章将介绍通过mysqldump、ibbackup、replication这些工具对数据进行备份的方法，同时对一些新兴的备份工具进行简单的介绍，如xtrabacup、LVM等。希望通过DBA的日常备份工作，能够保证在故障发生时，将数据丢失减少到最小。

# 8.1 备份与恢复概述

根据备份的方法不同可以将备份分为：

□Hot Backup（热备）

□Cold Backup（冷备）

□Warm Backup（温备）

Hot Backup是指数据库运行中直接备份，对正在运行的数据库操作没有任何的影响，这种方式在MySQL官方手册中称为Online Backup（在线备份）。Cold Backup是指备份操作是在数据库停止的情况下，这种备份最为简单，一般只需要复制相关的数据库物理文件即可，这种方式在MySQL官方手册中称为Offline Backup（离线备份）。Warm Backup备份同样是在数据库运行中进行的，但是会对当前数据库的操作有所影响，如加一个全局读锁以保证备份数据的一致性。

按照备份后文件的内容，备份又可以分为：

□逻辑备份

□裸文件备份

在MySQL数据库中，逻辑备份是指备份出的文件内容是可读的，一般是文本内容。内容一般是由一条条SQL语句，或者是表内实际数据组成。如mysqldump和SELECT*INTO OUTFILE的方法。这类方法的好处是可以观察导出文件的内容，一般适用于数据库的升级、迁移等工作。但其缺点是恢复所需要的时间往往较长。

裸文件备份是指复制数据库的物理文件，既可以是在数据库运行中的复制（如ibbackup、xtrabackup这类工具），也可以是在数据库停止运行时直接的数据文件复制。这类备份的恢复时间往往较逻辑备份短很多。

按照备份数据库的内容来分，备份又可以分为：

□完全备份

□增量备份

□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□（replay）□□□□□□□□point-in-time□□□□□□□□MySQL□□□□□（replication）□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□slave/standby□□□□□□□

□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□100GB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□LSN□□□□□□□□□□□LSN□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□xtrabackup□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□InnoDB□□□□□□□□□□□□MVCC□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□REPEATABLE READ□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□mysqldump□□□□□□□□□□□□□□□□□--single-transaction□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□--single-transaction□□□□□□□□□mysqldump□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□24×7□□□□□□□□□□□□□□□DAO（Database Admin Online）□□□□□□□□□□□□DBA□□□□□□□□□□Python□□□□□□Web□□□□□□□Django□□□□□□□□□DBA□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□DBA□

2008

## 8.2 　□□

由于InnoDB的表是独立表空间，每个表都有一个MySQL数据库的frm文件，用于存放表的相关信息，还有一个*.ibd文件，用于存放表的数据和索引。同时，MySQL还有一个全局配置文件my.cnf，用于配置数据库的各种参数。

作为DBA，需要很清楚地知道哪些数据需要备份，DBA不仅要保证数据的完整性和一致性，还需要保证备份数据的可恢复性。对数据进行备份是一项重要的工作，需要DBA有足够的耐心和细心，对每一个细节都不能放过。

数据备份是一项非常重要的工作，数据是一个企业最重要的资产之一，一旦数据丢失，将会给企业带来不可估量的损失。

常见的备份方式：

□逻辑备份：备份的是数据库的逻辑结构

□物理备份：备份的是数据库的物理文件，MySQL的数据文件和日志文件

□全量备份：对整个数据库进行一次完整的备份

□增量备份：备份自上次备份以来的SQL日志，恢复时需要重放日志

备份的注意事项：

□InnoDB存储引擎在备份时，需要保证数据的一致性，因此需要对数据进行加锁，同时还需要备份undo日志，以保证数据的一致性

□备份时需要考虑到数据库的性能，尽量在MySQL数据库空闲的时候进行备份，以减少对业务的影响

# 8.3 备份工具

## 8.3.1 mysqldump

mysqldump客户端最初由Igor Romanenko编写，它是数据库一种逻辑dump工具，经常被用于数据库备份。它生成一组MySQL语句，当重新执行这些MySQL语句时，将重建MySQL数据库（也可为Oracle、Microsoft SQL Server创建）。

mysqldump的命令为：

```
shell＞mysqldump[arguments]＞fle_name
```

备份服务器上所有的数据库，使用--all-databases选项：

```
shell＞mysqldump--all-databases＞dump.sql
```

备份服务器上多个数据库，使用--databases选项:

```
shell＞mysqldump--databases db1 db2 db3＞dump.sql
```

备份数据库test，对所有的表使用单个事务的方式：

```
[root@xen-server＞]#mysqldump--single-transaction test＞test_backup.sql
```

上面的命令是备份数据库test，备份过程中使用--single-transaction参数，将备份数据库的数据写入test_backup.sql文件。备份完成后，可以使用cat命令查看其中的内容：

```
[root@xen-server＞]#cattest_backup.sql

--MySQL dump 10.13 Distrib 5.5.1-m2,for unknown-linux-gnu(x86_64)

--
```

--Host:localhost Database:test

-------------------------------------------------------

--Server version 5.5.1-m2-log

……

--

--Table structure for table'a'

--

DROP TABLE IF EXISTS'a';

/*!40101 SET@saved_cs_client=@@character_set_client*/;

/*!40101 SET character_set_client=utf8*/;

CREATE TABLE'a'(

'b'int(11)NOT NULL DEFAULT'0',

PRIMARY KEY('b')

)ENGINE=InnoDB DEFAULT CHARSET=latin1;

/*!40101 SET character_set_client=@saved_cs_client*/;

--

--Dumping data for table'a'

--

LOCK TABLES'a'WRITE;

/*!40000 ALTER TABLE'a'DISABLE KEYS*/;

INSERT INTO'a'VALUES(1),(2),(4),(5);

/*!40000 ALTER TABLE'a'ENABLE KEYS*/;

UNLOCK TABLES;

--

--Table structure for table'z'

--

DROP TABLE IF EXISTS'z';

/*!40101 SET@saved_cs_client=@@character_set_client*/;

/*!40101 SET character_set_client=utf8*/;

CREATE TABLE'z'(

'a'int(11)DEFAULT NULL

)ENGINE=InnoDB DEFAULT CHARSET=latin1;

/*!40101 SET character_set_client=@saved_cs_client*/;

--

--Dumping data for table'z'

--

LOCK TABLES'z'WRITE;

/*!40000 ALTER TABLE'z'DISABLE KEYS*/;

INSERT INTO'z'VALUES(1),(1);

/*!40000 ALTER TABLE'z'ENABLE KEYS*/;

UNLOCK TABLES;

……

--Dump completed on 2010-08-03 13:36:17

从上面的文件内容可以看出，逻辑备份恢复是通过SQL语句的方法来完成的，这也是逻辑备份的特点，即通过MySQL存储的数据来重新生成文件，恢复的过程比较慢。上面显示的通过CREATE TABLE来创建表，通过INSERT等SQL来插入数据。

mysqldump有很多的选项，可以通过mysqldump--help命令来查看所有的参数，有些参数是默认启用的，如--lock-tables，可以通过-l来显示地启用和禁用某些参数。

□--single-transaction：在备份开始前，先执行START TRANSACTION命令，以此来获得备份的一致性，当前该参数只对InnoDB存储引擎有效。当启用该参数并进行备份时，确保没有其他任何的DDL语句执行，因为一致性读并不能隔离DDL操作。

□--lock-tables（-l）：在备份中，依次锁住每个架构下的所有表。一般用于MyISAM存储引擎，当只有只读访问时，保证备份的一致性。对于InnoDB存储引擎，不需要使用该参数，用--single-transaction即可。而且--lock-tables和--single-transaction是互斥（exclusive）的，不能同时使用。如果你的MySQL数据库中既有MyISAM存储引擎的表，又有InnoDB存储引擎的表，那么这时你的选择只有--lock-tables了。值得注意的是，参数--lock-tables并不能保证所有表的一致性，这是因为在备份时，架构可能会发生改变。

☐--lock-all-tables，-x：在开始转储之前为所有数据库中的所有表请求锁定。该选项自动关闭--lock-tables和单事务选项的取值。

☐--add-drop-database：在CREATE DATABASE前添加DROP DATABASE语句。默认在--all-databases或--databases选项所指的情况下，每个数据库名前添加CREATE DATABASE语句。下面为使用该选项的示例。

---

[root@xen-server～]#mysqldump--single-transaction--add-drop-database--databases test＞test_backup.sql

[root@xen-server～]#cat test_backup.sql

--MySQL dump 10.13 Distrib 5.5.1-m2,for unknown-linux-gnu(x86_64)

……

--

--Current Database:'test'

--

/*!40000 DROP DATABASE IF EXISTS'test'*/;

CREATE DATABASE/*!32312 IF NOT EXISTS*/'test'/*!40100 DEFAULT CHARACTER SET latin1*/;

USE'test';

……

---

☐--master-data[=value]：该选项将二进制的位置和文件名追加到输出文件中。这个选项将对replication产生影响。如果value的值为1，位置和文件名将被写到CHANGE MASTER语句中。如果value的值为2，在CHANGE MASTER前添加SQL注释。另外，如果value的值省略，value默认为1。下面为使用该选项的示例。

---

[root@xen-server～]#mysqldump--single-transaction--add-drop-database--master-data=1--databases test＞test_backup.sql

[root@xen-server～]#cat test_backup.sql

--MySQL dump 10.13 Distrib 5.5.1-m2,for unknown-linux-gnu(x86_64)

--

--Host:localhost Database:test

------------------------------------------------------

```
--Server version 5.5.1-m2-log

……

--

--Position to start replication or point-in-time recovery from

--

CHANGE MASTER TO MASTER_LOG_FILE='xen-server-bin.000006',MASTER_LOG_POS=8095;

……
```

---

## 当value为2时，可以得到一个注释掉的CHANGE MASTER语句，用于参考

---

```
[root@xen-server～]#mysqldump--single-transaction--add-drop-database--master-data=2--databases test＞test_backup.sql

[root@xen-server～]#cat test_backup.sql

--MySQL dump 10.13 Distrib 5.5.1-m2,for unknown-linux-gnu(x86_64)

--

--Host:localhost Database:test

--------------------------------------------------------

--Server version 5.5.1-m2-log

……

--

--Position to start replication or point-in-time recovery from

--

--

……
```

---

□--master-data选项会打开--lock-tables选项，使用本选项与--single-transaction选项时会打开--lock-all-tables选项。

□--events（-E）：备份相关的所有事件。

□--routines（-R）：备份存储过程及存储函数。

□--triggers：备份触发器。

□--hex-blob：将BINARY、VARBINARY、BLOG、BIT列类型备份为十六进制的格式。mysqldump的备份文件通常是文本文件，但是如果导出的数据中有二进制数据，就会出现乱码，可以使用--hex-blob选项使用十六进制格式来导出二进制数据。

---

```
[root@xen-server□]#mysqldump--single-transaction--add-drop-database--master-data=2--no-autocommit--databases test3□
test3_backup.sql

[root@xen-server□]#cat test3_backup.sql

--MySQL dump 10.13 Distrib 5.5.1-m2,for unknown-linux-gnu(x86_64)

--

--Host:localhost Database:test3

--------------------------------------------------------

--Server version 5.5.1-m2-log

……

LOCK TABLES'a'WRITE;

/*!40000 ALTER TABLE'a'DISABLE KEYS*/;

setautocommit=0;

INSERT INTO'a'VALUES(0x61000000000000000000);

/*!40000 ALTER TABLE'a'ENABLE KEYS*/;

UNLOCK TABLES;
```

---

从结果可以看出，0x61000000000000000000就是二进制数据的十六进制格式。

□--tab=path、-T path：产生TAB分割的数据文件。对于每个表，mysqldump创建一个包含CREATE TABLE语句的table_name.sql文件，和包含数据的tbl_name.txt文件。可以使用--fields-terminated-by=…、--fields-enclosed-by=…、--fields-optionally-enclosed-by=…、--fields-escaped-by=…、--lines-terminated-by=…来修饰分隔符和换行符的输出格式。

---

```
[root@xen-server test]#mysqldump--single-transaction--add-drop-database--tab="/usr/local/mysql/data/test"test
```

```
[root@xen-server test]#ls-lh

total 244K

-rw-rw----1 mysql mysql 8.4K Jul 21 16:02 a.frm

-rw-rw----1 mysql mysql 96K Jul 22 17:18 a.ibd

-rw-r--r--1 root root 1.3K Aug 3 15:36 a.sql

-rw-rw-rw-1 mysql mysql 8 Aug 3 15:36 a.txt

-rw-rw----1 mysql mysql 65 Jul 17 15:54 db.opt

-rw-rw----1 mysql mysql 8.4K Aug 2 17:22 z.frm

-rw-rw----1 mysql mysql 96K Aug 2 17:22 z.ibd

-rw-r--r--1 root root 1.3K Aug 3 15:36 z.sql

-rw-rw-rw-1 mysql mysql 4 Aug 3 15:36 z.txt

-----------

--Server version 5.5.1-m2-log

/*!40101 SET@OLD_CHARACTER_SET_CLIENT=@@CHARACTER_SET_CLIENT*/;

/*!40101 SET@OLD_CHARACTER_SET_RESULTS=@@CHARACTER_SET_RESULTS*/;

/*!40101 SET@OLD_COLLATION_CONNECTION=@@COLLATION_CONNECTION*/;

/*!40101 SET NAMES utf8*/;

/*!40103 SET@OLD_TIME_ZONE=@@TIME_ZONE*/;

/*!40103 SET TIME_ZONE='+00:00'*/;

/*!40101 SET@OLD_SQL_MODE=@@SQL_MODE,SQL_MODE=''*/;

/*!40111 SET@OLD_SQL_NOTES=@@SQL_NOTES,SQL_NOTES=0*/;

--

--Table structure for table'a'

--

DROP TABLE IF EXISTS'a';

/*!40101 SET@saved_cs_client=@@character_set_client*/;

/*!40101 SET character_set_client=utf8*/;

CREATE TABLE'a'(

'b'int(11)NOT NULL DEFAULT'0',

PRIMARY KEY('b')
```

)ENGINE=InnoDB DEFAULT CHARSET=latin1;

/*!40101 SET character_set_client=@saved_cs_client*/;

/*!40103 SET TIME_ZONE=@OLD_TIME_ZONE*/;

/*!40101 SET SQL_MODE=@OLD_SQL_MODE*/;

/*!40101 SET CHARACTER_SET_CLIENT=@OLD_CHARACTER_SET_CLIENT*/;

/*!40101 SET CHARACTER_SET_RESULTS=@OLD_CHARACTER_SET_RESULTS*/;

/*!40101 SET COLLATION_CONNECTION=@OLD_COLLATION_CONNECTION*/;

/*!40111 SET SQL_NOTES=@OLD_SQL_NOTES*/;

--Dump completed on 2010-08-03 15:36:56

[root@xen-server test]#cat a.txt

1

2

4

5

---

# 但是有时候，DBA要通过SELECT…INTO OUTFILE的方式来导出数据，而通过mysqldump也可以完成同样的功能，并且会更加安全，其语法如下

## 用--where='where_condition'或-w'where_condition'来导出符合条件的记录，如导出字段b的值大于a表中字段a的值大于2的

---

[root@xen-server bin]#mysqldump--single-transaction--where='b＞2'test a＞a.sql

[root@xen-server bin]#cat a.sql

--MySQL dump 10.13 Distrib 5.5.1-m2,for unknown-linux-gnu(x86_64)

--

--Host:localhost Database:test

-------------------------------------------------------

--Server version 5.5.1-m2-log

……

--

--Dumping data for table'a'

--

--WHERE:b￿2

LOCK TABLES'a'WRITE;

/*!40000 ALTER TABLE'a'DISABLE KEYS*/;

INSERT INTO'a'VALUES(4),(5);

/*!40000 ALTER TABLE'a'ENABLE KEYS*/;

UNLOCK TABLES;

/*!40103 SET TIME_ZONE=@OLD_TIME_ZONE*/;

……

# 8.3.2 SELECT...INTO OUTFILE

SELECT...INTO语句用于将数据从表中检索出来并存储到文件中。下面给出SELECT...INTO的语法。

---

```
SELECT[column 1],[column 2]...

INTO

OUTFILE'file_name'

[{FIELDS|COLUMNS}

[TERMINATED BY'string']

[[OPTIONALLY]ENCLOSED BY'char']

[ESCAPED BY'char']

]

[LINES

[STARTING BY'string']

[TERMINATED BY'string']

]

FROM TABLE WHERE......
```

---

其中FIELDS[TERMINATED BY'string']用来分隔每个字段的值，[[OPTIONALLY]ENCLOSED BY'char']用来将每个字段的值括起，[ESCAPED BY'char']为转义字符，[STARTING BY'string']为行起始符，TERMINATED BY'string'为行结束符。这里给出默认情况下FIELDS和LINES的设置，如下所示。

---

```
FIELDS TERMINATED BY'\t'ENCLOSED BY''ESCAPED BY'\\'

LINES TERMINATED BY'\n'STARTING BY''
```

---

file_name为导出文件的名字，此文件不能为已存在的文件，mysql:mysql用户是MySQL服务器的运行用户。

mysql：select*into outfile'/root/a.txt'from a;

ERROR 1(HY000):Can't create/write to file'/root/a.txt'(Errcode:13)

## 一般把数据导出到数据库目录下

[root@xen-server：]#mysql test-e"select*into outfile'/home/mysql/a.txt'fields terminated by','from a";

ERROR 1086(HY000)at line 1:File'/home/mysql/a.txt'already exists

## 下面用到SELECT INTO语句和表a的数据

mysql：select*into outfile'/home/mysql/a.txt'from a;

Query OK,3 rows affected(0.02 sec)

mysql：quit

Bye

[root@xen-server：]#cat/home/mysql/a.txt

1 a

2 b

3 c

## 数据行之间使用换行符分隔，使用TAB键来分隔字段。如果想使用其他字符比如"，"分隔字段，用FIELDS TERMINATED BY'string'选项来设置

[root@xen-server：]#mysql test-e"select*into outfile'/home/mysql/a.txt'fields terminated by','from a";

[root@xen-server：]#cat/home/mysql/a.txt

1,a

2,b

3,c

# 在Windows中文本的换行符是"\r\n"，所以在导出文本时可以指定LINES TERMINATED BY这个参数：

---

[root@xen-servermysql]#mysql test-e"select*into outfile'/home/mysql/a.txt'fields terminated by','lines terminated by'\r\n'from a";

[root@xen-servermysql]#od-c a.txt

0000000 1,a\r\n 2,b\r\n 3,c\r\n

0000017

---

# 8.3.3 □□□□□□□□

mysqldump□□□□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□

---

[root@xen-server□]#mysql-uroot-p□test_backup.sql

Enter password:

---

□□□□□□□□□□□□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□drop database test;

ERROR 1010(HY000):Error dropping database(can't rmdir'./test',errno:39)

---

□□□□□□□□□□□□□□SQL□□□□□□□□□□□□□□SOURCE□□□□□□□□□□□□□□□□□□□□□□□□□□□□

---

mysql□source/home/mysql/test_backup.sql;

Query OK,0 rows affected(0.00 sec)

Query OK,0 rows affected(0.00 sec)

……

Query OK,0 rows affected(0.00 sec)

Query OK,0 rows affected(0.00 sec)

---

□□mysqldump□□□□□□□□□□□□□□□□□□□□□□□□□□mysqldump□□□□□□□□□
□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□
□mysqldump□□□□□□□□□□□□□□□□□□□□□□□□□□□□□frm□□□□□□□□□□□□□□
□□□□□□□□□□□mysqldump□□□□□□□□□□□□□□

# 8.3.4　LOAD DATA INFILE

与使用mysqldump-tab方法等价的SELECT INTO OUTFILE语句，用来导入这个文件的等价语句是LOAD DATA INFILE。下面是LOAD DATA INFILE语句。

---

LOAD DATA INTO TABLE a IGNORE 1 LINES INFILE'/home/mysql/a.txt'

[REPLACE|IGNORE]

INTO TABLE tbl_name

[CHARACTER SET charset_name]

[{FIELDS|COLUMNS}

[TERMINATED BY'string']

[[OPTIONALLY]ENCLOSED BY'char']

[ESCAPED BY'char']

]

[LINES

[STARTING BY'string']

[TERMINATED BY'string']

]

[IGNORE number LINES]

[(col_name_or_user_var,...)]

[SET col_name=expr,...]

---

可以看到，两者之间LOAD DATA INFILE语句除了FILE文件名称以及导入的表不同之外，其他和SELECT INTO OUTFILE语句一样。而IGNORE number LINES选项可以忽略导入的前几行。下面是一个LOAD DATA INFILE的简单例子，并看看忽略文件的第一行数据。

---

mysql>load data infile'/home/mysql/a.txt'into table a;

Query OK,3 rows affected(0.00 sec)

Records:3 Deleted:0 Skipped:0 Warnings:0

---

# 不使用对InnoDB型外键约束需要在某个时间段关闭外键检查的可以这样处理。

mysql□SET@@foreign_key_checks=0;

Query OK,0 rows affected(0.00 sec)

mysql□LOAD DATA INFILE'/home/mysql/a.txt'INTO TABLE a;

Query OK,4 rows affected(0.00 sec)

Records:4 Deleted:0 Skipped:0 Warnings:0

mysql□SET@@foreign_key_checks=1;

Query OK,0 rows affected(0.00 sec)

# 还可以在加载数据的同时修改或者计算某些值，比如对a和b两个c列表示a、b列的和。

mysql□CREATE TABLE b(

-□a INT,

-□b INT,

-□c INT,

-□PRIMARY KEY(a)

-□)ENGINE=InnoDB;

Query OK,0 rows affected(0.01 sec)

mysql□LOAD DATA INFILE'/home/mysql/a.txt'

-□INTO TABLE b FIELDS TERMINATED BY',' (a,b)

-□SET c=a+b;

Query OK,4 rows affected(0.01 sec)

Records:4 Deleted:0 Skipped:0 Warnings:0

mysql□SELECT*FROM b;

+---+------+------+

|a|b|c|

+---+------+------+

|1|2|3|

```
|2|3|5|

|4|5|9|

|5|6|11|

+---+------+------+

4 rows in set(0.00 sec)
```

# 8.3.5　mysqlimport

mysqlimport是MySQL提供的一个命令行程序，它提供了LOAD DATA INFILE的命令接口，而且很多选项都与LOAD DATA INFILE命令相同。可以这样使用

---

shell＞mysqlimport[options]db_name textfile1[textfile2...]

---

和LOAD DATA INFILE不同，mysqlimport命令可以用来导入多张表。并且通过--user-thread参数并发地导入不同的文件。这里的并发是指并发导入多个文件，而不是指并发导入多张表到一张表。并且，mysqlimport的实现是通过命令行的方式。下面显示了通过mysqlimport并发导入2个文件

---

[root@xen-servermysql]#mysqlimport--use-threads=2 test/home/mysql/t.txt/home/mysql/s.txt

test.s:Records:5000000 Deleted:0 Skipped:0 Warnings:0

test.t:Records:5000000 Deleted:0 Skipped:0 Warnings:0

---

## 通过命令我们可以观察到，在MySQL数据库中有两个会话对应的两个进程在同时进行导入

---

mysql＞SHOW FULL PROCESSLIST\G;

***************************1.row***************************

Id:46

User:rep

Host:www.dao.com:1028

db:NULL

Command:Binlog Dump

Time:37651

State:Master has sent all binlog to slave;waiting for binlog to be updated

Info:NULL

***************************2.row***************************

Id:77

User:root

Host:localhost

db:test

Command:Query

Time:0

State:NULL

Info:show full processlist

***************************3.row***************************

Id:83

User:root

Host:localhost

db:test

Command:Query

Time:73

State:NULL

Info:LOAD DATA INFILE'/home/mysql/t.txt'INTO TABLE't'IGNORE 0 LINES

***************************4.row***************************

Id:84

User:root

Host:localhost

db:test

Command:Query

Time:73

State:NULL

Info:LOAD DATA INFILE'/home/mysql/s.txt'INTO TABLE's'IGNORE 0 LINES

4 rows in set(0.00 sec)

# 可以看到mysqlimport实际上是同时导入多个文件的，LOAD DTA INFILE语句则是其可以导

# 8.4 　备份和恢复的注意事项

对于二进制日志的一个最重要的应用是point-in-time恢复。二进制和MySQL支持的replication可以使得主服务器和从服务器之间备份的数据一致。因此，在配置文件中应开启二进制日志功能。

---

```
[mysqld]

log-bin=mysql-bin
```

---

在3.2.4节中已经介绍了，InnoDB存储引擎支持事务，因此必须要确保二进制日志中事务的提交和回滚是一致的。而使用InnoDB存储引擎必须将上述参数进行如下的修改：

---

```
[mysqld]

log-bin=mysql-bin

sync_binlog=1

innodb_support_xa=1
```

---

此外，在备份前还可以通过使用FLUSH LOGS命令来生成一个新的二进制日志文件，接着进行备份。

在恢复的过程中可以通过命令行工具mysqlbinlog，使用mysqlbinlog工具来恢复数据：

---

```
shell＞mysqlbinlog[options]log_fle...
```

---

例如，二进制binlog.0000001可以通过如下的命令来恢复：

---

```
shell＞mysqlbinlog binlog.0000001|mysql-uroot-p test
```

---

给定的文件一个接一个地被处理，并使用单一连接来处理所有的日志。下面就是一个例子：

---

```
shell：mysqlbinlog binlog.[0-10]*|mysql-u root-p test
```

---

## 另一种方式是mysqlbinlog把所有的日志写到一个单独的文件中，然后使用SOURCE命令处理这个文件的内容。下面就是一个例子：

---

```
shell：mysqlbinlog binlog.000001》/tmp/statements.sql

shell：mysqlbinlog binlog.000002》》/tmp/statements.sql

shell：mysql-u root-p-e"source/tmp/statements.sql"
```

---

## --start-position和--stop-position选项可以用于指定开始和结束的位置，这样你可以仅仅执行二进制日志的一部分。下面就是一个例子：

---

```
shell：mysqlbinlog--start-position=107856 binlog.0000001|mysql-uroot-p test
```

---

## --start-datetime和--stop-datetime选项可以用于指定二进制日志中事件的开始和结束的时间，效果类似--start-position和--stop-position选项。

# 8.5 备份

## 8.5.1 ibbackup

ibbackup是InnoDB存储引擎官方提供的热备工具，它可以同时备份MyISAM存储引擎和InnoDB存储引擎表。对于InnoDB存储引擎表，其备份工作原理如下。

1）记录备份开始时，InnoDB存储引擎重做日志文件的检查点LSN。

2）复制共享表空间文件以及独立表空间文件。

3）记录复制完表空间文件后，InnoDB存储引擎重做日志文件的检查点LSN。

4）复制在备份时产生的重做日志。

对于事务的数据库，如Microsoft SQL Server数据库和Oracle数据库，要实现数据库的备份，通常需要生成数据库快照，由数据库快照来实现备份。而InnoDB存储引擎并没有数据库快照的功能，因此ibbackup的备份是通过：

□ 在事务数据库上获取SQL操作。

□ 记录每个操作在重做日志文件中的对应项。

□ 回放每个操作对应的重做日志文件的内容。

□ 在不同的平台，ibbackup可以运行于Linux、Windows以及主流的UNIX系统平台。

ibbackup对InnoDB存储引擎表的备份的工作。

□ 在线备份，不锁表。

□ 增量备份，减少备份时间。

ibbackup有着良好的备份性能，其对备份InnoDB存储引擎表的速度取决于磁盘读取的速度。此外，用户可以使用开源工具Percona XtraBackup来对其进行替代。

XtraBackup的最大缺点是不支持ibbackup的增量备份功能,好在开源社区又给我们带来了希望,Percona推出的XtraBackup支持增量备份功能。

# 8.5.2 XtraBackup

XtraBackup是开源软件，由Percona公司开发维护，可用于对基于MySQL5.0或更新版本的XtraBackup以GPL v2协议发布，官方网站地址为
https://launchpad.net/percona-xtrabackup。

## xtrabackup的命令行语法格式

```
xtrabackup--backup|--prepare[OPTIONS]
```

## xtrabackup的命令行参数说明

(The defaults options should be given as the first argument)

--print-defaults Prints the program's argument list and exit.

--no-defaults Don't read the default options from any file.

--defaults-file=Read the default options from this file.

--defaults-extra-file=Read this file after the global options files have been read.

--target-dir=The destination directory for backups.

--backup Make a backup of a mysql instance.

--stats Calculate the statistic of the datadir(it is recommended you take mysqld offline).

--prepare Prepare a backup so you can start mysql server with your restore.

--export Create files to import to another database after it has been prepared.

--print-param Print the parameters of mysqld that you will need for a forcopyback.

--use-memory=This value is used instead of buffer_pool_size.

--suspend-at-end Creates a file called xtrabackup_suspended and waits until the user deletes that file at the end of the backup.

--throttle=(use with--backup)Limits the IO operations(pairs of reads and writes)per second to the values set here.

--log-stream outputs the contents of the xtrabackup_logfile to stdout.

--incremental-lsn=(use with--backup)Copy only.ibd pages newer than the specified LSN high:low.##ATTENTION##:checkpoint lsn*must*be used.Be Careful!

--incremental-basedir=(use with--backup)Copy only.ibd pages newer than the existing backup at the specified directory.

--incremental-dir=(use with--prepare)Apply.delta files and logfiles located in the specified directory.

--tables=name Regular Expression list of table names to be backed up.

--create-ib-logfile(NOT CURRENTLY IMPLEMENTED)will create ib_logfile*after a--prepare.

###If you want to create ib_logfile*only re-execute this

command using the same options.###

--datadir=name Path to the database root.

--tmpdir=name Path for temporary files.Several paths may be specified as a colon(:)separated string.

If you specify multiple paths they are used round-robin.

---

# 実際にバックアップとリストアを試してみる

---

#./xtrabackup--backup

./xtrabackup Ver alpha-0.2 for 5.0.75 unknown-linux-gnu(x86_64)

□□log scanned up to(0 1009910580)

Copying./ibdata1

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/ibdata1

...done

Copying./tpcc/stock.ibd

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/tpcc/stock.ibd

...done

Copying./tpcc/new_orders.ibd

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/tpcc/new_orders.ibd

...done

Copying./tpcc/history.ibd

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/tpcc/history.ibd

...done

Copying./tpcc/customer.ibd

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/tpcc/customer.ibd

□□log scanned up to(0 1010561109)

...done

Copying./tpcc/district.ibd

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/tpcc/district.ibd

...done

Copying./tpcc/item.ibd

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/tpcc/item.ibd

...done

Copying./tpcc/order_line.ibd

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/tpcc/order_line.ibd

□□log scanned up to(0 1012047066)

...done

Copying./tpcc/orders.ibd

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/tpcc/orders.ibd

...done

Copying./tpcc/warehouse.ibd

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp2/tpcc/warehouse.ibd

...done

□□log scanned up to(0 1014592707)

Stopping log copying thread..

Transaction log of lsn(0 1009910580)to(0 1014592707)was copied.

□□□□□□□□□□□□xtrabackup□□□□□□□□□□□□□□□□□□□□□□□□□□□□0 1009910580□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□copy□□□□□□□□□□□□□□Copying...to...□□□□□□□□□□□□□□□□□□□□□□0 1014592707□□

# 8.5.3  XtraBackup增量备份原理

MySQL本身提供的工具并不支持真正的增量备份，二进制日志恢复虽然是point-in-time（时间点）的恢复，但它其实是对上次的完全备份进行binlog的回放操作，而XtraBackup工具支持对于InnoDB存储引擎的增量备份，其工作原理如下。

1）首次完全备份，记录下此时检查点的LSN。

2）在进行增量备份时，比较表空间中每个页的LSN是否大于上次备份时的LSN，如果是，则备份该页，同时记录当前检查点的LSN。

因此，XtraBackup增量备份其备份和恢复的命令如下：

---

```
(full backup)

#./xtrabackup--backup--target-dir=/backup/base

...

(incremental backup)

#./xtrabackup--backup--target-dir=/backup/delta--incremental-basedir=/backup/base

...

(prepare)

#./xtrabackup--prepare--target-dir=/backup/base

...

(apply incremental backup)

#./xtrabackup--prepare--target-dir=/backup/base--incremental-dir=/backup/delta

...
```

---

首先进行一个完全备份，并指定目录/backup/base，然后进行增量备份，并指定目录/backup/delta。接着在恢复的过程中，先恢复完全备份，然后再应用增量备份进行恢复即可。最后需要注意的是：

---

```
#./xtrabackup--backup
```

./xtrabackup Ver beta-0.4 for 5.0.75 unknown-linux-gnu(x86_64)

□□log scanned up to(0 378161500)

...

The latest check point(for incremental):'0:377883685'□=====□□□□LSN

□□log scanned up to(0 379294296)

Stopping log copying thread..

Transaction log of lsn(0 377883685)to(0 379294296)was copied.

(must do--prepare before the each incremental backup)

#./xtrabackup--prepare

...

#./xtrabackup--backup--incremental=0:377883685

incremental backup from 0:377883685 is enabled.

./xtrabackup Ver beta-0.4 for 5.0.75 unknown-linux-gnu(x86_64)

□□log scanned up to(0 379708047)

Copying./ibdata1

to/home/kinoyasu/xtrabackup_work/mysql-5.0.75/innobase/xtrabackup/tmp_diff/ibdata1.delta

...done

...

The latest check point(for incremental):'0:379438233'□====□□□□□□□□□LSN

□□log scanned up to(0 380663549)

Stopping log copying thread..

Transaction log of lsn(0 379438233)to(0 380663549)was copied.

## 8.6　□□□□

MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□FreeBSD□UFS□□□□□Solaris□ZFS□□□□□GNU/Linux□□□□□□□Logical Volume Manager□LVM□□□□□□LVM□□□□□□□UFS□ZFS□□□□□□□□LVM□□□□

LVM□LINUX□□□□□□□□□□□□□□□□□□□□□LVM□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□LVM□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Volume Group□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Logical Volumes□□□□□□□□□□□□□□□□□□□□□□□LVM□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□LVM□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□8-1□□□□

Logical Volumes

lvcreate

Volume Group

vgcreate

Block Devices

pvcreate

图 8-1 LVM示意图

如8-2所示，此时已经建立了逻辑卷LV0。



图 8-2 建立逻辑卷示意图

下面vgdisplay命令将显示系统中卷组的信息。

[root@nh124-98□]#vgdisplay

---Volume group---

VG Name rep

System ID

Format lvm2

Metadata Areas 1

Metadata Sequence No 1873

VG Access read/write

VG Status resizable

MAX LV 0

Cur LV 3

Open LV 1

Max PV 0

Cur PV 1

Act PV 1

VG Size 260.77 GB

PE Size 4.00 MB

Total PE 66758

Alloc PE/Size 66560/260.00 GB

Free PE/Size 198/792.00 MB

VG UUID MQJiye-j4NN-LbZG-F3CQ-UdTU-fo9D-RRfXD5

---

vgdisplay命令可以看到上面我们创建的卷组rep的相关信息：260.77GB、访问权限为可读写read/write等。而lvdisplay命令则可以看到如下逻辑卷的信息：

---

[root@nh124-98~]#lvdisplay

---Logical volume---

LV Name/dev/rep/repdata

VG Name rep

LV UUID 7tOlDt-seKZ-ChpY-QMXC-WaFD-zXAl-MRbofK

LV Write Access read/write

LV snapshot status source of

/dev/rep/dho_datasnapshot100805143507[active]

/dev/rep/dho_datasnapshot100805163504[active]

LV Status available

#open 1

LV Size 100.00 GB

Current LE 25600

Segments 1

Allocation inherit

Read ahead sectors auto

-currently set to 256

Block device 253:0

---Logical volume---

LV Name/dev/rep/dho_datasnapshot100805143507

VG Name rep

LV UUID fSSXzh-IBnZ-aZIn-eP03-b7pk-CPjN-5xUktE

LV Write Access read only

LV snapshot status active destination for/dev/rep/repdata

LV Status available

#open 0

LV Size 100.00 GB

Current LE 25600

COW-table size 80.00 GB

COW-table LE 20480

Allocated to snapshot 0.13%

Snapshot chunk size 4.00 KB

Segments 1

Allocation inherit

Read ahead sectors auto

-currently set to 256

Block device 253:1

---Logical volume---

LV Name/dev/rep/dho_datasnapshot100805163504

VG Name rep

LV UUID 3B9NP1-qWVG-pfJY-Bdgm-DIdD-dUMu-s2L6qJ

LV Write Access read only

LV snapshot status active destination for/dev/rep/repdata

LV Status available

#open 0

```
LV Size 100.00 GB

Current LE 25600

COW-table size 80.00 GB

COW-table LE 20480

Allocated to snapshot 0.02%

Snapshot chunk size 4.00 KB

Segments 1

Allocation inherit

Read ahead sectors auto

-currently set to 256

Block device 253:4
```

---

从上面的命令可以看出3个逻辑卷的信息，卷组rep还剩余空间。其中逻辑卷大小为100GB的/dev/rep/repdata，是需要进行快照备份的逻辑卷，这里不再赘述。

LVM快照其实就是Copy-on-write（写时复制）的一种实现，也就是快照创建的时候，并不会将源逻辑卷的数据（meta data）真正地复制一份到快照逻辑卷中，而只是复制了源逻辑卷的一些组织结构信息，这个时候占用的空间是很少的；当源逻辑卷的数据发生改变的时候，会将源逻辑卷中被改写的数据，在改写之前先复制到快照逻辑卷中，这样就保证了快照逻辑卷中的数据和创建快照时刻源逻辑卷中的数据始终保持一致。

图8-3描述了LVM快照的机制，分区B是源逻辑卷中的数据，在创建快照后如果发生改变，则分区A、C、D的内容不变，只有分区B的数据会发生改变。

图 8-3 LVM快照工作

使用lvcreate时,可以通过增加参数--permission r选项来创建只读的快照卷。

```
[root@nh119-215 data]#lvcreate--size 100G--snapshot--permission r-n datasnapshot/dev/rep/repdata
```

Logical volume"datasnapshot"created

可以通过逻辑卷的查看命令lvdisplay查看到所创建的快照的COW-table size即所创建的快照的空间大小，以及Allocated to snapshot即已经使用的快照的空间比率。

[root@nh124-98~]#lvdisplay

……

---Logical volume---

LV Name/dev/rep/dho_datasnapshot100805163504

VG Name rep

LV UUID 3B9NP1-qWVG-pfJY-Bdgm-DIdD-dUMu-s2L6qJ

LV Write Access read only

LV snapshot status active destination for/dev/rep/repdata

LV Status available

#open 0

LV Size 100.00 GB

Current LE 25600

COW-table size 80.00 GB

COW-table LE 20480

Allocated to snapshot 0.04%

Snapshot chunk size 4.00 KB

Segments 1

Allocation inherit

Read ahead sectors auto

-currently set to 256

Block device 253:4

从上面可以看出已经使用了0.04%的快照的空间，一旦创建的快照的空间被占满之后，创建的快照就会被自动删除掉，相应的快照恢复功能也就失效了。

用LVM创建的InnoDB□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

由于InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 8.7 □□

## 8.7.1 □□□□□□□□

□□□replication□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□replication□□□□□□□□□□3□□□□□

1□□□□□□□master□□□□□□□□□□□□□□□□□□□□□□binlog□□□□

2□□□□□□□slave□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□relay log□□□□

3□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□8-4□□□□

主服务器

数据
更改

二进制
日志

日志传送

I/Q 线程

从服务器

SQL 线程

写入

读取

中继日志

# 图 8-4 MySQL复制进程的应用状态

从图可见，2个系统进程分别是I/O进程和应用中继日志的进程，它们都处于等待接收主库同步数据的状态，SQL进程中继日志应用（在MySQL4.0以后复制进程被拆分成1个数据同步进程和应用中继日志进程）。而且用户进程是执行SHOW命令的进程。以上进程在运行过程中的状态值是不断变化的，具体的状态值含义可参考官方文档。

<hr>

mysql＞SHOW FULL PROCESSLIST\G;

***************************1.row***************************

Id:1

User:system user

Host:

db:NULL

Command:Connect

Time:6501

State:Waiting for master to send event

Info:NULL

***************************2.row***************************

Id:2

User:system user

Host:

db:NULL

Command:Connect

Time:0

State:Has read all relay log;waiting for the slave I/O thread to update it

Info:NULL

***************************3.row***************************

Id:206

User:root

Host:localhost

db:NULL

Command:Query

Time:0

State:NULL

Info:SHOW FULL PROCESSLIST

3 rows in set(0.00 sec)

---

其中，ID为1的线程是I/O线程，它负责与主库通信接收其二进制日志；ID为2的线程是SQL线程，它负责读取并执行接收到的二进制日志内容。上述输出反映出I/O线程是以复制（replication）用户的身份与主库通信，主库会为每个从库创建一个专门的连接线程。

---

mysql＞SHOW FULL PROCESSLIST\G;

……

***************************65.row***************************

Id:26541

User:rep

Host:192.168.190.98:39549

db:NULL

Command:Binlog Dump

Time:6857

State:Has sent all binlog to slave;waiting for binlog to be updated

Info:NULL

……

---

为了进一步查看MySQL主从复制的运行状态，我们可以在主库和从库上分别执行命令SHOW SLAVE STATUS和SHOW MASTER STATUS进行查看。

---

mysql＞SHOW SLAVE STATUS\G;

***************************1.row***************************

Slave_IO_State:Waiting for master to send event

Master_Host:192.168.190.10

```
Master_User:rep

Master_Port:3306

Connect_Retry:60

Master_Log_File:mysql-bin.000007

Read_Master_Log_Pos:555176471

Relay_Log_File:gamedb-relay-bin.000048

Relay_Log_Pos:224355889

Relay_Master_Log_File:mysql-bin.000007

Slave_IO_Running:Yes

Slave_SQL_Running:Yes

Replicate_Do_DB:

Replicate_Ignore_DB:

Replicate_Do_Table:

Replicate_Ignore_Table:

Replicate_Wild_Do_Table:

Replicate_Wild_Ignore_Table:mysql.%,DBA.%

Last_Errno:0

Last_Error:

Skip_Counter:0

Exec_Master_Log_Pos:555176471

Relay_Log_Space:224356045

Until_Condition:None

Until_Log_File:

Until_Log_Pos:0

Master_SSL_Allowed:No

Master_SSL_CA_File:

Master_SSL_CA_Path:

Master_SSL_Cert:

Master_SSL_Cipher:

Master_SSL_Key:
```

Seconds_Behind_Master:0

Master_SSL_Verify_Server_Cert:No

Last_IO_Errno:0

Last_IO_Error:

Last_SQL_Errno:0

Last_SQL_Error:

1 row in set(0.00 sec)

在用SHOW SLAVE STATUS查看服务器的状态时,包括的参数如表8-1所示。

表 8-1  SHOW SLAVE STATUS 的主要变量

| 变　量 | 说　明 |
|---|---|
| Slave_IO_State | 显示当前 IO 线程的状态，上述状态显示的是等待主服务发送二进制日志 |
| Master_Log_File | 显示当前同步的主服务器的二进制日志，上述显示当前同步的是主服务器的 mysql-bin.000007 |
| Read_Master_Log_Pos | 显示当前同步到主服务器上二进制日志的偏移量位置，单位是字节。上述的示例显示当前同步到 mysql-bin.000007 的 555176471 偏移量位置，即已经同步了 mysql-bin.000007 这个二进制日志中 529MB（555176471/1024/1024）的内容 |
| Relay_Master_Log_File | 当前中继日志同步的二进制日志 |
| Relay_Log_File | 显示当前写入的中继日志 |
| Relay_Log_Pos | 显示当前执行到中继日志的偏移量位置 |
| Slave_IO_Running | 从服务器中 IO 线程的运行状态，YES 表示运行正常 |
| Slave_SQL_Running | 从服务器中 SQL 线程的运行状态，YES 表示运行正常 |
| Exec_Master_Log_Pos | 表示同步到主服务器的二进制日志偏移量的位置。（Read_Master_Log_Pos - Exec_Master_Log_Pos）可以表示当前 SQL 线程运行的延时，单位是字节。上述例子显示当前主从服务器是完全同步的 |

通过SHOW MASTER STATUS语句可以查看主服务器中二进制日志的状态，如下所示。

```
mysql〉SHOW MASTER STATUS\G;
```

```
**************************1.row**************************

File:mysql-bin.000007

Position:606181078

Binlog_Do_DB:

Binlog_Ignore_DB:

1 row in set(0.01 sec)
```

可以看到，主数据库服务器的位置是606181078，与从数据库服务器上通过监控看到的Read_Master_Log_Pos是一致的，说明I/O线程是正常的。

主从复制是通过MySQL自带的二进制日志同步来实现的，通过两个线程，分别是I/O线程和SQL线程共同完成，只要其中有一个线程出现异常，主从复制就会出现故障，有可能导致主从数据不一致的情况发生。

## 8.7.2 复制+分区的架构

下面是一些将复制和分区相结合的场景，供读者参考。

❑ 在分区后，每个MySQL数据库服务器可以再搭建多台从服务器，以满足读取的需要，同时也起到备份的作用。

❑ 因为对数据进行了分区，所以每台数据库服务器上的数据量得到了降低。在高并发的情况下，通过DNS、Round-Robin、Linux的LVS技术来实现负载均衡。

❑ 当某台数据库服务器产生问题时，可以方便地切换到其他数据库服务器。

❑ 将每台数据库服务器的复制再延迟几个小时，以防止人为的误操作，如主服务器上的数据库被删除。

当然，上述方案也存在缺点。例如，由于分区后，数据库的架构变得很复杂，维护也相应变得困难。同时，读和写的压力均分散到了每台服务器，因此写的压力反而会变大，这又变相地增加了数据库的负载。此外，由于数据库在不同的服务器上，所以如果不使用分布式事务，数据库的一致性将很难得到保证。

在结束本章之前，我们要对复制（replication）技术再说几句。很多开发人员会有这样的困惑，以为存在了复制，数据库就可高枕无忧了。然而事实并非如此，复制技术本质上只是实现了point-in-time的恢复，复制+分区的架构见图8-5所示。

图 8-5 主从+快照服务器工作

通常，只有具有相应权限的用户才可以在数据库中执行写入操作，其他用户只具有读权限。但是在某些突发情况下，也可能需要将所有用户（包括那些具有写权限的用户）设置为只读状态，例如数据库维护、备份或在服务器之间进行数据迁移。

可以使用如下命令使数据库处于read-only状态，该命令需要在配置文件中进行设置，保存配置文件之后重新启动数据库。

---

[mysqld]

read-only

---

执行了read-only命令之后，数据库处于只读状态，只有具有SUPER权限的用户才可以进行写入操作，其他用户只可以进行读取操作。

---

mysql>INSERT INTO z SELECT 2;

ERROR 1290(HY000):The MySQL server is running with the--read-only option so it cannot execute this statement

---

# 8.8 　小结

本章对数据库的备份恢复进行了介绍，MySQL数据库的不同存储引擎备份恢复的方法也略有不同，InnoDB存储引擎本章重点介绍了mysqldump、xtrabackup等工具。此外还对InnoDB存储引擎的快照复制进行了描述，希望读者在阅读完本章后，在任何突发的情况下都能从容地恢复数据库，保证整个数据库服务与MySQL数据库的可用性。

# 第9章 性能优化

一个系统的性能瓶颈可能来自各个方面，可能是网络层面的，可能是InnoDB存储引擎层面的，可能是操作系统层面的，可能是程序代码层面的，可能是InnoDB存储引擎参数层面的等，本章将主要介绍关于InnoDB存储引擎的性能优化。

□ 选择合适的CPU

□ 内存的重要性

□ 硬盘对数据库性能的影响

□ 合理地设置RAID

□ 操作系统的选择也很重要

□ 不同的文件系统对数据库的影响

□ 选择合适的基准测试工具

# 9.1　合理地选择CPU

在选择CPU时，用户应该选择多核还是高主频的CPU？用户应该了解应用类型，是OLTP（Online Transaction Processing，在线事务处理）还是OLAP（Online Analytical Processing，在线分析处理）应用，以此来选择合适的CPU。如果是OLAP应用，建议选择多核的处理器；同样，选择支持并行的SQL查询也是有益的。对于OLTP应用，比如大部分的网页应用、电子商务应用、Blog等类型的网站，建议选择的是OLAP类型的数据库应用。

InnoDB存储引擎一般都应用于OLTP的数据库应用，这种应用的特点如下所示。

□ 用户操作的并发量大

□ 事务处理的时间一般比较短

□ 查询的语句较为简单，一般都走索引

□ 复杂的查询较少

可以发现，上述这些OLTP的数据库应用的特点决定了CPU的选择，多核的CPU显然是最好的选择。因此对于多核的CPU，数据库应用都是OLTP，多核处理是最有益的。而对于OLAP，CPU单核处理器。而对于OLTP，IO设备都是应用关心的，因此可以说，这种处理器对于IO操作都有益。

从购买或升级数据库服务器时，使用多核的CPU一般都是64位的处理器。因为64位的处理器能访问的内存大于64位的CPU能访问到的。因此对于4核的CPU来说，市面上的Intel、AMD已经可以生产8核的CPU，能够访问大量的内存，因此现在的128核的CPU已经被大量应用于数据库服务器。

从InnoDB存储引擎的设计架构上来看，其并不能完全发挥多核的特性。因为master thread的工作主要集中在主线程中，因此并不能发挥多核性能的特性。不过随着InnoDB1.0的推出，这个情况已经得到了改善，多核CPU可以被更好的利用起来。而到了InnoDB 1.2版本时，将purge操作单独放到一个线程，从而master thread中所需要运行的线程减少了CPU。此外，在InnoDB存储引擎，从1.1版本开始，可以支持多个CPU的读写操作。用户可以设置innodb_read_io_threads和innodb_write_io_threads，增加IO的处理线程，从而充分利用多核CPU的处理性能。

□□□□MySQL□□□□□□□□□□□□SQL□□□□□□□□□□□□□□CPU□□□□□□□□□□□□□□CPU□□□□□
OLTP□□□□□□□□□□□□□□□□□□□□□□□□□OLTP□□□□□□□□□□□□□□□□□□□□□□□CPU□□□□□
CPU□□□□□□□□□□□□□□□□□□□□□□□□□

## 9.2　数据库性能

对于关系数据库来说，数据库最大的性能问题就是磁盘访问。因为InnoDB引擎支持数据缓存，缓存数据越多，越有助于提升数据库的整体性能。在InnoDB Buffer Pool中，数据量的大小对于数据库性能有着明显的影响。Percona公司的CTO Vadim曾经就内存对于数据库整体性能的影响进行过测试，如图9-1所示。

sysbench oltp,80mln rows（18GB data）

# 例 9-1 如何合理地设置InnoDB缓冲池的容量大小

假设某个应用的数据量大概在18GB，我们可以将缓冲池的容量大小设置为2GB、4GB、6GB、8GB、10GB、12GB、14GB、16GB、18GB、20GB、22GB，然后用sysbench等工具进行性能压测，直到系统的整体性能TPS（Transaction Per Second）趋于稳定，比如设置为20GB、22GB时系统的整体性能趋于稳定，此时对应的缓冲池容量大小就可以作为合理的缓冲池容量大小，这个值需要根据具体的应用场景进行合理的调整。

当然，缓冲池的容量也不是"越大"越好，缓冲池容量大小的设置还要受限于服务器本身的内存容量大小，其默认值为64MB（字节）。

另外，还可以通过缓冲池命中率这个指标查看缓冲池的设置是否合理。有一个经验值可供参考，一般情况下，InnoDB缓冲池的命中率不应该小于99%，如下：

---

```
mysql＞SHOW GLOBAL STAUTS LIKE'innodb%read%'\G;

***************************1.row***************************

Variable_name:Innodb_buffer_pool_read_ahead

Value:0

***************************2.row***************************

Variable_name:Innodb_buffer_pool_read_ahead_evicted

Value:0

***************************3.row***************************

Variable_name:Innodb_buffer_pool_read_requests

Value:167051313

***************************4.row***************************

Variable_name:Innodb_buffer_pool_reads

Value:129236

***************************5.row***************************

Variable_name:Innodb_data_pending_reads

Value:0

***************************6.row***************************

Variable_name:Innodb_data_read

Value:2135642112
```

***************************7.row***************************

Variable_name:Innodb_data_reads

Value:130309

***************************8.row***************************

Variable_name:Innodb_pages_read

Value:130215

***************************9.row***************************

Variable_name:Innodb_rows_read

Value:17651085

9 rows in set(0.00 sec)

各变量的含义如表9-1所示。

表 9-1  当前服务器的状态参数

| 参　数 | 说　明 |
|---|---|
| Innodb_buffer_pool_reads | 表示从物理磁盘读取页的次数 |
| Innodb_buffer_pool_read_ahead | 预读的次数 |
| Innodb_buffer_pool_read_ahead_evicted | 预读的页，但是没有被读取就从缓冲池中被替换的页的数量，一般用来判断预读的效率 |
| Innodb_buffer_pool_read_requests | 从缓冲池中读取页的次数 |
| Innodb_data_read | 总共读入的字节数 |
| Innodb_data_reads | 发起读取请求的次数，每次读取可能需要读取多个页 |

通过以上参数，可以得出以下两个指标。

缓冲池命中率

$$= \frac{\text{Innodb\_buffer\_pool\_read\_requests}}{\left(\text{Innodb\_buffer\_pool\_read\_requests} + \text{Innodb\_buffer\_pool\_read\_ahead} + \text{Innodb\_buffer\_pool\_reads}\right)}$$

$$平均每次读取的字节数 = \frac{\text{Innodb\_data\_read}}{\text{Innodb\_data\_reads}}$$

根据上述数据可以计算出缓冲池命中率=167 051 313/（167 051 313+129 236+0）=99.92%。

由此可以看出，缓冲池命中率非常高。如果缓冲池命中率低于百分之九十，则应该考虑扩充内存，增加缓冲池空间。这样可以使更多的数据在内存中，降低磁盘读取的操作。

## 9.3 □□□□□□□□□□□□

### 9.3.1 □□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□SAS□SATA□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□2.5□□SAS□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□3ms□□□□15 000RPM（rotate per minute）□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□RAID□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 9.3.2 □□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Flash Memory□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□EMC□□□□□□□□□□□□□□□□□□□□□TB□□□□□□□□□□□□□□Oracle□□□□□□□□□□□□□□□□□□Exadata□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□sector□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□erase□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□128KB□□256KB□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□9-2□□□□□□□□□□□□□□□□□□□□□□4□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Intel X-25M □□□□□□□10□□□□□□□□□

图 9-2 闪存控制芯片的体系结构

这么高的性能。但是，或许令大家感到意外的是，目前市场上固态硬盘的成本正在快速大幅度下降，而闪存的出现为我们存储层次结构又增加了一层。由于闪存是固态的，因此它在诸如随机读取方面的性能比传统机械硬盘好的多。同时，较低廉的价格也使得企业可以使用多块固态硬盘组成 I/O 性能非常强劲的磁盘阵列。当然，使用闪存作为存储也有缺点，如读写速度的不均匀性，这也对数据库的存储提出了新的挑战。通常来说，内存的随机读取延时小于 0.1ms，如图 9-3 所示。而闪存是机械硬盘，随机读取延时则显示出不同的情况。

Random Access Time（ms）



图 9-3 各种存储介质在随机访问时间上的差异显示

由于固态硬盘的高性能，InnoDB 存储引擎中的几个参数，如 innodb_io_capacity，就需要根据固态硬盘进行调整。该参数表示磁盘的 IOPS 能力，默认值可能已经无法发挥固态硬盘的高性能了。此外，InnoSQL 和 InnoDB1.2 版本都根据固态硬盘的特性进行了更深层次的优化，以充分利用固态硬盘的高性能特性。

另一个比较好的方法是，InnoSQL 提供 L2 Cache 的功能，将传统机械硬盘上的部分热点数据存储在固态硬盘上，从而提高数据库整体的读取性能。可以这样认为，内存是一级缓冲，固态硬盘是二级缓冲，而机械硬盘为最终存储层。由于固态硬盘

这里不作展开，仅列出部分实现了类似Cache功能的软件产品：Facebook Flash Cache、bcache、还有商业产品。实际上，最新版本的InnoDB中也已经实现了类似

# 9.4 磁盘阵列（RAID

## 9.4.1 RAID概述

RAID（Redundant Array of Independent Disks）即独立磁盘冗余阵列，磁盘阵列是由很多价格较便宜的磁盘，组合成一个容量巨大的磁盘组，利用个别磁盘提供数据所产生加成效果提升整个磁盘系统效能。利用这项技术，将数据切割成许多区段，分别存放在各个硬盘上。RAID技术主要包含冗余技术和条带技术。磁盘阵列还能利用同位检查的观念，在数组中任意一个硬盘故障时，仍可读出数据，在数据重构时，将数据经计算后重新置入新硬盘中。

RAID主要优点：

□ 传输速率高。

□ 可靠性高。

□ 磁盘容量大。

根据实际情况选择适当的RAID级别可以满足用户对存储系统可用性、性能和容量的要求。常用的RAID级别有以下几种：RAID 0、RAID 1、RAID 5、RAID 10、RAID 50等。

RAID 0连续以位或字节为单位分割数据，并行读/写于多个磁盘上，因此具有很高的数据传输率，但它没有数据冗余，如图9-4所示。RAID 0只是单纯地提高性能，并没有为数据的可靠性提供保证，而且其中的一个磁盘失效将影响到所有数据。因此，RAID 0不能应用于数据安全性要求高的场合。由于不带校验的RAID 0成本低，要求至少两个磁盘，我们通常只是为了提高数据读/写速度的需要而采用RAID 0，其在所有的级别中，速度是最快的。但是RAID 0既没有冗余功能，也不具备容错能力，如果一个磁盘（物理）损坏，所有数据都会丢失，危险程度与JBOD相当。例如，假设单个磁盘读取速度为50MB/s，两块磁盘RAID 0为96MB/s，四块磁盘RAID 0为130MB/s而非150MB/s。

図　9-4　RAID 0原理

RAID 1的原理是把N个磁盘分成两组（如图9-5所示），一组是数据盘，另一组是数据盘的镜像盘，也就是说数据同时写入两组磁盘，两组磁盘的内容完全一样。RAID 1拥有和单个磁盘一样的存储容量，实际可用的存储空间只是总存储空间的一半，可靠性是所有磁盘阵列中最高的。由于是数据镜像，RAID 1的数据传输速率和单个磁盘基本一样。RAID 1被磁盘阵列中成本最高的一种，因此RAID 1常用在对可靠性要求较高的场合。

图 9-5 RAID 1示意

RAID 5是一种存储性能、数据安全和存储成本兼顾的存储解决方案。它使用的是Disk Striping（硬盘分区）技术。RAID 5至少需要三块硬盘，RAID 5不是对存储的数据进行备份，而是把数据和相对应的奇偶校验信息存储到组成RAID 5的各个磁盘上，并且奇偶校验信息和相对应的数据分别存储于不同的磁盘上。当RAID 5的一个磁盘数据发生损坏后，利用剩下的数据和相应的奇偶校验信息去恢复被损坏的数据。RAID 5可以理解为是RAID 0和RAID 1的折中方案。RAID 5可以为系统提供数据安全保障，但保障程度要比镜像低而磁盘空间利用率要比镜像高。RAID 5具有和RAID 0相近似的数据读取速度，只是多了一个奇偶校验信息，写入数据的速度相当的慢，若使用Write Back可以让性能改善不少。同时由于多个数据对应一个奇偶校验信息，RAID 5的磁盘空间利用率要比RAID 1高，存储成本相对较低。RAID 5的结构如图9-6所示。

RAID 5



图 9-6　RAID 5原理

RAID 10和RAID 01。RAID 10是将镜像和条带进行组合的一种磁盘阵列，先用RAID 0将磁盘进行分区，再用RAID 1将分区进行镜像。RAID 10是将镜像和条带进行组合的一种磁盘阵列，先用RAID 0将磁盘进行分区，再用RAID 01来实现。RAID 10是将镜像和条带进行组合的一种磁盘阵列，RAID 01先用磁盘进行镜像，再用RAID 1将磁盘进行分区，先镜像再用RAID 0分区。RAID 01和RAID 10这两种组合方式在实际应用中都很常见，它们之间的主要区别在于RAID 01先进行镜像再用RAID 10先进行RAID 0分区再进行镜像。RAID 10先进行RAID 0分区再RAID 1镜像，它们之间的主要区别在于组合的先后顺序不同，其具体实现方式如图RAID 10和RAID 01分别如图9-7所示。

图 9-7    RAID 10、RAID 01对比

RAID 50：RAID 50在本质上是将多组磁盘阵列先建立成各自独立的RAID 0硬盘阵列，然后再将这几组磁盘阵列组合成RAID 5，这样既提高了读取速度，也增加了安全性，同时也可以提供比较好的数据恢复能力。RAID 5的实现较为……

RAID 50



图 9-8 RAID 50示意

与此类似的还有人们常说的RAID 10，它是先做镜像也就是做RAID 1再RAID 0。而本例讲的则是将每
一个磁盘按照一定的大小分成若干个条带（strip），然后再对不同磁盘上的条带做镜像。这样做的

即在保持高速状态下所得到的流量，这时磁盘的IO效率为20%；另一种为磁盘连续读写时IO效率，达到90%。

# 9.4.2　RAID Write Back□□

RAID Write Back□□□□□RAID□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□sync_binlog□□1□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□Write Back□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□RAID□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□RAID□□□□□□□□□□□□BBU□Battery Backup Unit□□□□□□□□□□□□□□Write Back□□□□□□□□□□□□□□□□□□□□□□□□□□□RAID□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□Write Back□□□□□□RAID□□□□□□□□□□Write Through□Write Through□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□Write Back□□□RAID□□□□□□□□Write Through□□□□□□□□□□□□□□Write Back□□□□□RAID□□□□□□□□□□□□□□□□□□□□□□□RAID□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Write Back□□□□□□□□□□□Write Through□

□□□□□□□□□□□□□□□□□□□□□□□□Write Back□□□□□□□□□□□□□□□□□□□□□Write Back□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Write Back□

□□□□□□□□20W□□□□□□□Write Back□Write Through□□□□□□□

---

```
mysql□CREATE TABLE t(a CHAR(2))Engine=InnoDB;

Query OK,0 rows affected(0.00 sec)

mysql□DELIMITER//

mysql□

mysql□CREATE PROCEDURE p()

-□BEGIN

-□DECLARE v INT;
```

```
-	SET v=0;

-	WHILE v	200000 DO

-	INSERTINTO t VALUES('aa');

-	SET v=v+1;

-	END WHILE;

-	END

-	//

Query OK,0 rows affected(0.12 sec)

mysql	DELIMITER;
```

接着向这个表中插入t插入20W条记录，我们分别测试了Write Back、Write Through对性能带来的影响。测试结果如表9-2所示。

表 9-2　Write Back 和 Write Through 的性能对比测试结果

| RAID 卡设置 | 时　间 |
| --- | --- |
| Write Back | 43 秒 |
| Write Through | 31 分钟 |
| Write Through with innodb_flush_log_at_trx_commit=0 | 68 秒 |

测试结果显示，性能上的差距还是非常大的。通过存储过程CALL P，向表中插入同样的20W条记录，但是所需的时间相差很大。Write Back模式下仅仅需要43秒，而Write Through模式则需要31分钟，达到了40倍的差距。

我们可以对Write Through模式进行优化，例如将参数innodb_flush_log_at_trx_commit设置为0，再次执行存储过程P，这时所需的时间为68秒。虽然优化后的性能提升非常明显，但是相对于原来的性能还是有所下降。如果是在master数据库上，那么强烈不建议这样的设置，因为可能会影响事务的数据完整性和安全性。

# 9.4.3  RAID□□□□

□RAID□□□□□□□□□□□□□□□□□□□□□□□BIOS□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□RAID□□□□□□□□□□□□□□LSI□□□□□□□□RAID□□□□□□□MegaCLI□□□□□□□□□

MegaCLI□□□□□□□□□□□□□□□□Windows□□□□□□□□□GUI□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□MegaCLI□□□□□Windows□□□□□□□□□□MegaCLI.exe□

□□MegaCLI□□RAID□□□□□□

---

[root@xen-server□]#/opt/MegaRAID/MegaCli/MegaCli64-AdpAllInfo-a0

Adapter#0

============================================================================

Versions

================

Product Name:MegaRAID SAS 8708ELP

Serial No:P012233608

FW Package Build:9.0.1-0030

……

HW Configuration

================

SAS Address:500605b000d1e180

BBU:Present

Alarm:Present

NVRAM:Present

Serial Debugger:Present

Memory:Present

Flash:Present

Memory Size:256MB

TPM:Absent

Default Settings

=================

Phy Polarity:0

PhyPolaritySplit:240

Background Rate:30

Stripe Size:64kB

Flush Time:4 seconds

Write Policy:WB

Read Policy:None

Cache When BBU Bad:Disabled

Cached IO:No

SMART Mode:Mode 6

Alarm Disable:Yes

Coercion Mode:1GB

ZCR Config:Unknown

Dirty LED Shows Drive Activity:No

BIOS Continue on Error:No

Spin Down Mode:None

Allowed Device Type:SAS/SATA Mix

Allow Mix in Enclosure:Yes

Allow HDD SAS/SATA Mix in VD:Yes

Allow SSD SAS/SATA Mix in VD:No

Allow HDD/SSD Mix in VD:No

Allow SATA in Cluster:No

Max Chained Enclosures:3

Disable Ctrl-R:Yes

Enable Web BIOS:Yes

Direct PD Mapping:No

BIOS Enumerate VDs:Yes

Restore Hot Spare on Insertion:No

Expose Enclosure Devices:Yes

Maintain PD Fail History:Yes

Disable Puncturing:No

Zero Based Enclosure Enumeration:No

PreBoot CLI Enabled:Yes

LED Show Drive Activity:No

Cluster Disable:Yes

SAS Disable:No

Auto Detect BackPlane Enable:SGPIO/i2c SEP

Use FDE Only:No

Enable Led Header:No

Delay during POST:0

---

从上面长长的输出信息中可以获得很多有价值的信息，比如RAID卡的型号和缓存大小。本例中RAID卡型号为MegaRAID SAS 8708ELP，缓存大小为256MB。以及缓存策略等信息，本例中Write Policy为WB（Write Back）。

MegaCLI查看当前物理磁盘信息的命令如下：

---

[root@xen-server ]#/opt/MegaRAID/MegaCli/MegaCli64-PDList-aALL

Adapter#0

Enclosure Device ID:252

Slot Number:0

Device Id:8

Sequence Number:2

Media Error Count:0

Other Error Count:0

Predictive Failure Count:0

Last Predictive Failure Event Seq Number:0

PD Type:SAS

Raw Size:279.396 GB[0x22ecb25c Sectors]

Non Coerced Size:278.896 GB[0x22dcb25c Sectors]

Coerced Size:278.464 GB[0x22cee000 Sectors]

Firmware state:Online

SAS Address(0):0x5000c5000f363b55

SAS Address(1):0x0

Connected Port Number:0(path0)

Inquiry Data:SEAGATE ST3300655SS 00023LM5MGZZ

FDE Capable:Not Capable

FDE Enable:Disable

Secured:Unsecured

Locked:Unlocked

Foreign State:None

Device Speed:Unknown

Link Speed:Unknown

Media Type:Hard Disk Device

……

---

从上面的输出可以知道，这是SEAGATE ST3300655SS的硬盘，具体的硬盘信息可以在这个网站上查到http://discountechnology.com/Seagate-ST3300655SS-SAS-Hard-Drive，它的尺寸大小为3.5英寸，转速为15000转，它的Cache为16MB，总之这是一个普通的3.5英寸硬盘，接口类型为SAS4.0的硬盘。

可以看到，两个逻辑磁盘的缓存策略都是Write Back的。

---

[root@xen-server ~]#/opt/MegaRAID/MegaCli/MegaCli64-LDGetProp-Cache-LALL-aALL

Adapter 0-VD 0(target id:0):Cache Policy:WriteBack,ReadAheadNone,Direct,No Write Cache if bad BBU

Adapter 0-VD 1(target id:1):Cache Policy:WriteBack,ReadAheadNone,Direct,No Write Cache if bad BBU

Exit Code:0x00

---

所有逻辑磁盘（虚拟磁盘、RAID）的Write Back（前提是服务器有BBU。若没有，请不要设置成Write Back，可能丢失数据，因为缓存在RAID卡上，而非硬盘上）、读策略和磁盘缓存策略：

以下两条命令分别为开启和关闭写策略：

---

#/opt/MegaRAID/MegaCli/MegaCli64-LDSetPropWB-LALL-aALL

#/opt/MegaRAID/MegaCli/MegaCli64-LDSetPropWT-LALL-aALL

---

以上命令开启所有逻辑磁盘（RAID）的写回。将Write Back（写回）Write Through（直写）。注意，要想将所有Write Through（直写）Write Back（写回），请将上面的命令互换。

# 9.5 操作系统的选择

Linux是MySQL数据库服务器的主流操作系统，大多数用户都选择使用Linux作为操作系统平台。□□□□□□□□□□□□□□□□□□□□□□□□□Linux□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

除了Linux操作系统，FreeBSD也是一个很好的选择。□□□□□□□□□□□□FreeBSD的MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□FreeBSD的MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

Solaris□□□□□□□□□□□□□□□□□□□SPARC□□□□□□□□□□□□□□□□□□□□X86□□□□□Solaris□□□□□□□□□□□□□□□□□□□□□□□□ZFS□□□□□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□Open Solaris□

Windows□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Windows□□□□MySQL□□□□□□□□□□□□□□□□□□□□□Linux□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Windows□□□□□□□□□□□□□□□□□□□Linux□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

4G□□□□□□□□□□□□□□□□□□□□□□□□□□□□8G□□□□□□□□□□□□□□□□□□□□4G□□□□□□□□□□□□□□□□64□□□□□□□□□□□□□□□□□□□□□□□□□□64□□□□□□□□□□□□□64□□□□□□□□□□□□□□64□□□□□□□□□□□□□□□□□□□□□□□□□□32□□MySQL□□□□□□□□64□□□□□□□□□□□□□□□□□□64□□□□□□□□□□□□□□□□□

## 9.6 □□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□Windows□□□□□NTFS□Solaris□□□□ZFS□□□□□Linux□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□EXT3□□□□□ReiserFS□□□□□EXT4□□□□□XFS□

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□XFS□□□□□□"□□"□□□□□□□□□□□□□□□□□□□□□□□□□□□□□EXT3□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□EXT3□□□□□□□□□□□□□□□□DBA□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□□□□□□□□□□□□□□□□□□DBA□□□□□□□□□□ZFS□□□□□□□□□□□□□□□□□□□□□□□□□□□□LVM□□□□□□□□□□□□□□□□□□□□□□□mount□□□□□□□□□□□□□□□□□□□□□□□□□□□□

# 9.7　常见数据库性能测试工具

数据库上线之前需要进行比较详细的性能测试，尤其是使用MySQL数据库，因为它的很多存储引擎都有各自的应用场景，所以选择合适的测试工具及方法尤其重要。下面介绍两款测试工具：sysbench和mysql-tpcc。

## 9.7.1　sysbench

sysbench是一个模块化、跨平台、多线程基准测试工具，主要用于评估测试各种不同系统参数下的数据库负载情况。它主要包括以下几种方式的测试：

□ CPU性能

□ 磁盘IO性能

□ 调度程序性能

□ 内存分配及传输速度

□ POSIX线程性能

□ 数据库性能（OLTP基准测试）

sysbench主要用于OLTP基准测试MySQL、PostgreSQL、Oracle等，目前sysbench主要支持Linux操作系统，另外，可以使用sysbench来测试Windows操作系统下的Microsoft SQL Server的数据库负载。

sysbench可以通过官网地址http://sysbench.sourceforge.net来获取最新版本。很多版本的sysbench已经收录到各个版本的操作系统软件库中。如果你使用的是Linux操作系统，那么可以使用RED HAT等软件库中已经编译好的sysbench版本，以节约安装时间。

sysbench的功能较多，执行下面的命令可以看到它支持的详细参数：

---

```
[root@xen-server ~]#sysbench

Missing required command argument.

Usage:
```

sysbench[general-options]....-test=[]test-name[][test-options]...command

General options:

--num-threads=N number of threads to use[1]

--max-requests=N limit for total number of requests[10000]

--max-time=N limit for total execution time in seconds[0]

--thread-stack-size=SIZE size of stack per thread[32K]

--init-rng=[on|off]initialize random number generator[off]

--test=STRING test to run

--debug=[on|off]print more debugging info[off]

--validate=[on|off]perform validation checks where possible[off]

--help=[on|off]print help and exit

--version=[on|off]print version and exit

Compiled-in tests:

fileio-File I/O test

cpu-CPU performance test

memory-Memory functions speed test

threads-Threads subsystem performance test

mutex-Mutex performance test

oltp-OLTP test

Commands:prepare run cleanup help version

See'sysbench--test=[]name[]help'for a list of options for each test.

---

# 由于InnoDB存储引擎是事务型的,所以我们以测试OLTP为例。但是在这之前,我们先以fileio和oltp为例,介绍怎么使用sysbench这个工具进行测试的。

---

[root@xen-server[]]#sysbench--test=fileio help

sysbench 0.4.10:multi-threaded system evaluation benchmark

fileio options:

--file-num=N number of files to create[128]

--file-block-size=N block size to use in all IO operations[16384]

--file-total-size=SIZE total size of files to create[2G]

--file-test-mode=STRING test mode{seqwr,seqrewr,seqrd,rndrd,rndwr,rndrw}

--file-io-mode=STRING file operations mode{sync,async,fastmmap,slowmmap}[sync]

--file-extra-flags=STRING additional flags to use on opening files{sync,dsync,direct}[]

--file-fsync-freq=N do fsync()after this number of requests(0-don't use fsync())[100]

--file-fsync-all=[on|off]do fsync()after each write operation[off]

--file-fsync-end=[on|off]do fsync()at the end of test[on]

--file-fsync-mode=STRING which method to use for synchronization{fsync,fdatasync}[fsync]

--file-merged-requests=N merge at most this number of IO requests if possible(0-don't merge)[0]

--file-rw-ratio=N reads/writes ratio for combined test[1.5]

---

各个参数的含义如下：

☐--file-num：表示生成测试文件的数量，默认为128。

☐--file-block-size：测试时所使用文件块的大小，如果想磁盘对InnoDB存储引擎进行测试，那么可以将它设置为16384，即InnoDB存储引擎页的大小。默认值为16384。

☐--file-total-size：创建测试文件的总大小，默认为2GB。

☐--file-test-mode：文件测试模式,包含seqwr（顺序写）、seqrewr（顺序读写）、seqrd（顺序读）、rndrd（随机读）、rndwr（随机写）、rndrw（随机读写）。

☐--file-io-mode：文件操作的模式，例如同步、异步、MMAP（map映射）等模式，默认为同步模式。

☐--file-extra-flags：打开文件时的选项，与API相关的参数。

☐--file-fsync-freq：执行fsync函数的频率。fsync主要是同步磁盘文件，因为可能有系统和磁盘缓冲的关系。

☐--file-fsync-all：每执行完一次写操作，就执行一次fsync，默认为off。

☐--file-fsync-end：在测试结束时执行fsync函数，默认为on。

☐--file-fsync-mode：文件同步的方式。该工具会调用API来进行文件同步，默认为fdatasync。可供选择的参数还有：fdatasync，同步fsync。

☐--file-rw-ratio：读写的比例，默认为2:1。

sysbench中fileio的执行分为prepare、run、cleanup三个阶段。prepare阶段主要用于准备测试文件，run是执行测试，cleanup是清理测试文件。下面的例子中，以16个文件总大小2GB进行fileio测试。

---

```
[root@xen-server ssd]#sysbench--test=fileio--file-num=16--file-total-size=2G prepare

sysbench 0.4.10:multi-threaded system evaluation benchmark

16 files,131072Kb each,2048Mb total

Creating files for the test...
```

---

## 上面的命令会创建总共16个文件，合计总大小为2GB，每个文件的大小平均为128MB。

---

```
[root@xen-server ssd]#ls-lh

total 2G

-rw-------1 root root 128M Aug 12 10:42 test_file.0

-rw-------1 root root 128M Aug 12 10:42 test_file.1

-rw-------1 root root 128M Aug 12 10:42 test_file.10

-rw-------1 root root 128M Aug 12 10:42 test_file.11

-rw-------1 root root 128M Aug 12 10:42 test_file.12

-rw-------1 root root 128M Aug 12 10:42 test_file.13

-rw-------1 root root 128M Aug 12 10:42 test_file.14

-rw-------1 root root 128M Aug 12 10:42 test_file.15

-rw-------1 root root 128M Aug 12 10:42 test_file.2

-rw-------1 root root 128M Aug 12 10:42 test_file.3

-rw-------1 root root 128M Aug 12 10:42 test_file.4

-rw-------1 root root 128M Aug 12 10:42 test_file.5

-rw-------1 root root 128M Aug 12 10:42 test_file.6
```

-rw-------1 root root 128M Aug 12 10:42 test_file.7

-rw-------1 root root 128M Aug 12 10:42 test_file.8

-rw-------1 root root 128M Aug 12 10:42 test_file.9

## 测试随机读。以下代码表示使用直接方式、16个线程进行测试,测试时间

[root@xen-server ssd]#sysbench--test=fileio--file-total-size=2G--file-test-mode=rndrd--max-time=180--max-requests=100000000--num-threads=16--init-rng=on--file-num=16--file-extra-flags=direct--file-fsync-freq=0--file-block-size=16384 run

## 为180秒,随机请求数最大为100 000 000个。但是,180秒内不可能完成这么多次的随机读操作,所以测试会在180秒时停止。

[root@xen-server ssd]#sysbench--test=fileio--file-total-size=2G--file-test-mode=rndrd--max-time=180--max-requests=100000000--num-threads=16--init-rng=on--file-num=16--file-extra-flags=direct--file-fsync-freq=0--file-block-size=16384 run

sysbench 0.4.10:multi-threaded system evaluation benchmark

Running the test with following options:

Number of threads:16

Initializing random number generator from timer.

Extra file open flags:16384

16 files,128Mb each

2Gb total file size

Block size 16Kb

Number of random requests for random IO:100000000

Read/Write ratio for combined random IO test:1.50

Calling fsync()at the end of test,Enabled.

Using synchronous I/O mode

Doing random read test

Threads started!

Time limit exceeded,exiting...

(last message repeated 15 times)

Done.

Operations performed:619908 Read,0 Write,0 Other=619908 Total

Read 9.459Gb Written 0b Total transferred 9.459Gb(53.81Mb/sec)

3443.85 Requests/sec executed

Test execution summary:

total time:180.0044s

total number of events:619908

total time taken by event execution:2878.0750

per-request statistics:

min:0.42ms

avg:4.64ms

max:27.30ms

approx.95 percentile:8.13ms

Threads fairness:

events(avg/stddev):38744.2500/102.69

execution time(avg/stddev):179.8797/0.00

---

# 测试结果表明数据吞吐为53.81MB/s，每秒执行IOPS为3443.85，其中事物执行最小消耗为每次消耗时间为，平均消耗时间为，最大消耗时间为。测试完成后，会生成测试数据文件，需通过cleanup进行清理，具体操作命令如下所示：

---

[root@xen-server ssd]#sysbench--test=fileio--file-num=16--file-total-size=2G cleanup

sysbench 0.4.10:multi-threaded system evaluation benchmark

Removing test files...

---

# 如果需要对磁盘进行反复测试，可通过简单的脚本来实现自动化测试，将测试结果进行输出并保存，以备后续使用。下面简单介绍一个磁盘自动化测试脚本的编写，具体如下所示：

---

#!/bin/sh

set-u

```
set-x

set-e

for size in 8G 64G;do

for mode in seqrd seqrw rndrd rndwr rndrw;do

for blksize in 4096 16384;do

sysbench--test=fileio--file-num=64--file-total-size=$size prepare

for threads in 1 4 8 16 32;do

echo"======testing$blksize in$threads threads"

echo PARAMS$size$mode$threads$blksize sysbench-size-$size-mode-$mode-threads-$threads-blksz-$blksize

for i in 1 2 3;do

sysbench--test=fileio--file-total-size=$size--file-test-mode=$mode\

--max-time=180--max-requests=100000000--num-threads=$threads--init-rng=on\

--file-num=64--file-extra-flags=direct--file-fsync-freq=0--file-block-size=$blksize run\

|tee-a sysbench-size-$size-mode-$mode-threads-$threads-blksz-$blksize 2  1

done

done

sysbench--test=fileio--file-total-size=$size cleanup

done

done

done
```

---

# 对于MySQL来说，除了OLTP测试，还有fileio测试。同样有prepare、run、cleanup。对于prepare，就会预先准备好数据文件，比如一个sbtest的文件。对于sbtest，sysbench可以假象成一个巨大的InnoDB表，比如下面的8000W数据。

---

sysbench 0.4.10:multi-threaded system evaluation benchmark

Creating table'sbtest'...

Creating 80000000 records in table'sbtest'...

---

# 用如下的语句进行数据库的oltp综合测试

---

sysbench--test=oltp--oltp-table-size=80000000--oltp-read-only=off--init-rng=on--num-threads=16--max-requests=0--oltp-dist-type=uniform--max-time=3600--mysql-user=root--mysql-socket=/tmp/mysql.sock--db-driver=mysql run〉res

---

# 该命令执行完后会生成一个res文件，在res文件中的部分内容如下：

---

sysbench 0.4.10:multi-threaded system evaluation benchmark

WARNING:Preparing of"BEGIN"is unsupported,using emulation

(last message repeated 15 times)

Running the test with following options:

Number of threads:16

Initializing random number generator from timer.

Doing OLTP test.

Running mixed OLTP test

Using Uniform distribution

Using"BEGIN"for starting transactions

Using auto_inc on the id column

Threads started!

Time limit exceeded,exiting...

(last message repeated 15 times)

Done.

OLTP test statistics:

queries performed:

read:6043324

write:2158330

other:863332

total:9064986

transactions:431666(119.90 per sec.)

deadlocks:0(0.00 per sec.)

read/write requests:8201654(2278.07 per sec.)

other operations:863332(239.80 per sec.)

Test execution summary:

total time:3600.2672s

total number of events:431666

total time taken by event execution:57598.5965

per-request statistics:

min:6.84ms

avg:133.43ms

max:7155.61ms

approx.95 percentile:325.84ms

Threads fairness:

events(avg/stddev):26979.1250/64.14

execution time(avg/stddev):3599.9123/0.06.

---

这里可以看到，在一个小时的连续测试中，transactions（即每秒钟处理的事务数，也就是TPS）的平均值达到了119.9tps。通过上面的这些数据，可以看出sysbench的OLTP测试对于TPS的测量很精确。而sysbench本身也具备了比较完善的功能，不仅可以进行数据库的

## 9.7.2　mysql-tpcc

TPC（Transaction Processing Performance Council）是一系列事务处理和数据库基准测试的标准规范，其中的TPC-C是TPC组织提供的一个基准测试标准，用来衡量OLTP系统的性能和可伸缩性，它是数据库的联机交易处理系统，TPC-C模拟的OLTP数据库应用系统。

TPC-C以3NF实现了一个较为复杂的订单处理事务模拟，主要涉及下订单、查询订单状态、发货、付款、检查仓库库存状态几个业务，其中包括一系列的数据库表，例如库存表有10万条记录、仓库表一般有3000条记录、订单表有100 000条记录等。这些数据库表记录着测试中的客户订单、库存数量等信息，它们之间的关系如图9中，本例中的关系如图9-9所示。

**stock**
- s_i_id INT(11)
- s_w_id SMALLINT(6)
- s_quantity SMALLINT(6)
- s_dist_01 CHAR(24)
- s_dist_02 CHAR(24)
- s_dist_03 CHAR(24)
- s_dist_04 CHAR(24)
- s_dist_05 CHAR(24)
- s_dist_06 CHAR(24)
- s_dist_07 CHAR(24)
- s_dist_08 CHAR(24)
- s_dist_09 CHAR(24)
- s_dist_10 CHAR(24)
- s_ytd DECIMAL(8,0)
- s_order_cnt SMALLINT(6)
- s_remote_cnt SMALLINT(6)
- s_data VARCHAR(50)
- Indexes

**warehouse**
- w_id SMALLINT(6)
- w_name VARCHAR(10)
- w_street_1 VARCHAR(20)
- w_street_2 VARCHAR(20)
- w_city VARCHAR(20)
- w_state CHAR(2)
- w_zip CHAR(9)
- w_tax DECIMAL(4,2)
- w_ytd DECIMAL(12,2)
- Indexes

**district**
- d_id TINYINT(4)
- d_w_id SMALLINT(6)
- d_name VARCHAR(10)
- d_street_1 VARCHAR(20)
- d_street_2 VARCHAR(20)
- d_city VARCHAR(20)
- d_state CHAR(2)
- d_zip CHAR(9)
- d_tax DECIMAL(4,2)
- d_ytd DECIMAL(12,2)
- d_next_o_id INT(11)
- Indexes

**history**
- h_c_id INT(11)
- h_c_d_id TINYINT(4)
- h_c_w_id SMALLINT(6)
- h_d_id TINYINT(4)
- h_w_id SMALLINT(6)
- h_date DATETIME
- h_amount DECIMAL(6,2)
- h_data VARCHAR(24)
- Indexes

**customer**
- c_id INT(11)
- c_d_id TINYINT(4)
- c_w_id SMALLINT(6)
- c_first VARCHAR(16)
- c_middle CHAR(2)
- c_last VARCHAR(16)
- c_street_1 VARCHAR(20)
- c_street_2 VARCHAR(20)
- c_city VARCHAR(20)
- c_state CHAR(2)
- c_zip CHAR(9)
- c_phone CHAR(16)
- c_since DATETIME
- c_credit CHAR(2)
- c_credit_lim BIGINT(20)
- c_discount DECIMAL(4,2)
- c_balance DECIMAL(12,2)
- c_ytd_payment DECIMAL(12,2)
- c_payment_cnt SMALLINT(6)
- c_delivery_cnt SMALLINT(6)
- c_data TEXT
- Indexes

**order_line**
- ol_o_id INT(11)
- ol_d_id TINYINT(4)
- ol_w_id SMALLINT(6)
- ol_number TINYINT(4)
- ol_i_id INT(11)
- ol_supply_w_id SMALLINT(6)
- ol_delivery_d DATETIME
- ol_quantity TINYINT(4)
- ol_amount DECIMAL(6,2)
- ol_dist_info CHAR(24)
- Indexes

**item**
- i_id INT(11)
- i_im_id INT(11)
- i_name VARCHAR(24)
- i_price DECIMAL(5,2)
- i_data VARCHAR(50)
- Indexes

**orders**
- o_id INT(11)
- o_d_id TINYINT(4)
- o_w_id SMALLINT(6)
- o_c_id INT(11)
- o_entry_d DATETIME
- o_carrier_id TINYINT(4)
- o_ol_cnt TINYINT(4)
- o_all_local TINYINT(4)
- Indexes

**new_orders**
- no_o_id INT(11)
- no_d_id TINYINT(4)
- no_w_id SMALLINT(6)
- Indexes

# 例 9-9　TPC-C基准程序

TPC-C的结果集以tpmC（tpm是transaction per minute，而C指TPC中C的基准程序）为单位，这代表了系统的吞吐量。

tpcc-mysql是由著名TPC-C基准程序的评测工具，可以通过TPC-C官方网站下载，地址是https://code.launchpad.net/～percona-dev/perconatools/tpcc-mysql，但是tpcc-mysql只能运行在Linux操作系统下，可以通过地址（Windows下载地址http://code.google.com/p/david-mysql-tools/downloads/list）下载Windows下的tpcc-mysql。

tpcc-mysql包括两个比较重要的部分：

□tpcc_load：根据仓库数（这里有9个仓库）来装载基准数据。

□tpcc_start：根据指定的参数进行TPC-C测试。

tpcc_load的提示信息如下所示：

---

```
[root@xen-server～]#tpcc_load

************************************

***###easy###TPC-C Data Loader***

************************************

usage:tpcc_load[server][DB][user][pass][warehouse]

OR

tpcc_load[server][DB][user][pass][warehouse][part][min_wh][max_wh]

*[part]:1=ITEMS 2=WAREHOUSE 3=CUSTOMER 4=ORDERS
```

---

其中各个参数含义如下：

□server：数据库MySQL服务器IP。

□DB：装载数据的数据库。

□user：MySQL数据库用户。

□pass□MySQL□□□□

□warehouse□□□□□□□□□□□

□□□tpcc_load□□□□□100□□□□□□□□tpcc□□□□□□□

---

[root@xen-server tpcc-mysql]#mysql tpcc□create_table.sql

[root@xen-server tpcc-mysql]#mysql tpcc□add_fkey_idx.sql

[root@xen-server tpcc-mysql]#tpcc_load 127.0.0.1 tpcc2 root xxxxxx 100

*************************************

***###easy###TPC-C Data Loader***

*************************************

□Parameters□

[server]:127.0.0.1

[DBname]:tpcc2

[user]:root

[pass]:

[warehouse]:100

TPCC Data Load Started...

Loading Item

...............................................5000

...............................................10000

...............................................15000

……□□□

...DATA LOADING COMPLETED SUCCESSFULLY.

---

## tpcc_start□□□□□□□□□□□□

---

[root@xen-server□]#tpcc_start

*************************************

***###easy###TPC-C Load Generator***

```
************************************

usage:tpcc_start[server][DB][user][pass][warehouse][connection][rampup][measure]
```

其中最后三个参数的含义如下：

□connection：虚拟的用户终端数量

□rampup：系统在多长时间后进入测试压力稳定状态（单位为秒）

□measure：测试的持续时间（单位为秒）

例如，tpcc_start以16个用户终端模拟压力，并在10分钟后系统压力稳定并持续20分钟，命令如下：

```
[root@xen-server]#tpcc_start 127.0.0.1 tpcc root xxxxxx 100 16 600 1200

************************************
***###easy###TPC-C Load Generator***
************************************
<Parameters>
[server]:127.0.0.1
[DBname]:tpcc
[user]:root
[pass]:xxxxxx
[warehouse]:100
[connection]:16
[rampup]:600(sec.)
[measure]:1200(sec.)
……
```

另外，测试过程中的状态及结果输出也十分容易理解，命令如下：

```
RAMP-UP TIME.(1 sec.)

MEASURING START.
```

10,624(0):0.4,624(0):0.2,62(0):0.2,63(0):0.6,62(0):0.8

20,990(0):0.2,988(0):0.2,98(0):0.2,99(0):0.4,98(0):0.6

30,1435(0):0.2,1436(0):0.2,144(0):0.2,143(0):0.2,144(0):0.4

40,1736(0):0.2,1739(0):0.2,174(0):0.2,174(0):0.2,174(0):0.4

50,2041(0):0.2,2044(0):0.2,204(0):0.2,204(0):0.2,207(0):0.2

60,2195(0):0.2,2193(0):0.2,220(0):0.2,221(0):0.2,218(0):0.2

70,2332(0):0.2,2335(0):0.2,233(0):0.2,232(0):0.2,234(0):0.2

80,2408(0):0.2,2401(0):0.2,241(0):0.2,239(0):0.2,241(0):0.2

90,2473(0):0.2,2476(0):0.2,247(0):0.2,250(0):0.2,248(0):0.2

100,2350(0):0.2,2347(0):0.2,235(0):0.2,233(0):0.2,235(0):0.2

……

在看这个图10的TPC-C测试之前，先来了解下TPC-C测试。它是5个事务类型：New Order、Payment、Order-Status、Delivery、Stock-Level，其中的New Order代表了TPC-C测试的核心事务类型，因此New Order Per 10 Second就作为了该测试的吞吐量指标。接下来的内容就是测试关于不同版本的InnoDB存储引擎10的测试，如图9-10所示。

TPCC



图 9-10  New Order Per 10  Second

在tpcc_load运行完毕后查看tpmC值，这里New Order Per 10 Second的稳定值，由于New Order Per 10 Second为每十秒数据，乘以6即得到每分钟的tpmC。

……

〖Constraint Check〗(all must be[OK])

[transaction percentage]

Payment:43.48%(⬚=43.0%)[OK]

Order-Status:4.35%(⬚=4.0%)[OK]

Delivery:4.35%(⬚=4.0%)[OK]

Stock-Level:4.35%(⬚=4.0%)[OK]

[response time(at least 90%passed)]

New-Order:99.72%[OK]

Payment:99.95%[OK]

Order-Status:99.93%[OK]

Delivery:100.00%[OK]

Stock-Level:100.00%[OK]

【TpmC】

7949.942 TpmC

# 9.8　小结

本章的前半部分介绍了InnoDB存储引擎的性能，主要从CPU、内存、磁盘这三部分来讲解。对于RAID，重点需要了解的内容是RAID的几种类型，以及写入时InnoDB存储引擎需要注意的问题，以及如何选择一个合适的磁盘阵列卡。对于Linux操作系统，讲解了文件系统的选择。sysbench、tpcc-mysql这两个工具是非常好的测试工具，尤其是后者可以帮助DBA测试某台服务器数据库MySQL所能承受的最大负载压力。

# 第10章 InnoDB存储引擎中的锁和多版本并发控制

InnoDB存储引擎在内部对于并发的控制就像它对大多数其他事情的处理方式一样：先理解Why，再看What。本章将对一系列的问题进行分析，包括各种锁的类型以及什么时候使用它们、多版本并发控制的工作原理以及InnoDB存储引擎是如何处理那些大多数数据库系统都会遇到的经典并发问题的。

# 10.1 认识InnoDB存储引擎的构成

InnoDB采用多层的架构，向上对接MySQL服务器，向下负责数据在MySQL中的落地。MySQL的总体架构<sup>[1]</sup>及其在MySQL中的位置如图10-1所示。

# MySQL Community Server 5.1.49

Select Platform:

| Source Code ▼ | Select |

| | | | |
|---|---|---|---|
| **SuSE Linux Enterprise Server ver. 11 (Architecture Independent), RPM Package** (MySQL-community-5.1.49-1.sles11.src.rpm) | 5.1.49 | 22.0M | Download MD5: 714c5f8bf4b1816bb88951649d9298aa |
| **Red Hat & Oracle Enterprise Linux 5 (Architecture Independent), RPM Package** (MySQL-community-5.1.49-1.rhel5.src.rpm) | 5.1.49 | 22.0M | Download MD5: 8c386345d6374be174033f3a17d49a0b |
| **SuSE Linux Enterprise Server 10 (Architecture Independent), RPM Package** (MySQL-community-5.1.49-1.sles10.src.rpm) | 5.1.49 | 22.0M | Download MD5: ca2ed7f15fae60331f40b0083847fe39 |
| **Generic Linux (glibc 2.3) (Architecture Independent), RPM Package** (MySQL-5.1.49-1.glibc23.src.rpm) | 5.1.49 | 22.0M | Download MD5: 40cb7b0399b3a174e9b9e06281f1036f |
| **SuSE Linux Enterprise Server 9 (Architecture Independent), RPM Package** (MySQL-community-5.1.49-1.sles9.src.rpm) | 5.1.49 | 22.0M | Download MD5: 9e171c70c47d6792b0856021bbd12353 |
| **Red Hat & Oracle Enterprise Linux 4 (Architecture Independent), RPM Package** (MySQL-community-5.1.49-1.rhel4.src.rpm) | 5.1.49 | 22.0M | Download MD5: 897d2e9ffbe9109072fe22bcafa91788 |
| **Red Hat Enterprise Linux 3 (Architecture Independent), RPM Package** (MySQL-community-5.1.49-1.rhel3.src.rpm) | 5.1.49 | 22.0M | Download MD5: bc855486bd2b4d5d029d7b6bba4d4361 |
| **Generic Linux (Architecture Independent),** | 5.1.49 | 22.6M | Download |

## 图 10-1 MySQL版本选择

接着在下面选择对应的操作系统类型，这里选择的是Generic Linux，这就是通用的
MySQL版本。在Download下载区域中可以看到GA和开发版两种下载链接，这里下载的
MySQL就是MySQL 5.5.5版本，milestone版本就是非GA版本，也就是开发版本，表
示这个版本还不够成熟，不推荐在生产环境中使用。注意区分www和dev两个站点，
http://dev.mysql.com/downloads/mysql这个下载链接就是下载MySQL
源代码的链接，如图10-2所示。

# MySQL Community Server 5.5.5 m3

Select Platform:

Source Code ▼ | Select

| | | | |
|---|---|---|---|
| **Linux - Generic 2.6 (Architecture Independent), RPM Package** | 5.5.5 | 21.0M | Download |
| (MySQL-5.5.5_m3-1.linux2.6.src.rpm) | | MD5: ca368fb09817420b2a0f3d325f2acc34 | |
| **SuSE Linux Enterprise Server ver. 11 (Architecture Independent), RPM Package** | 5.5.5 | 21.0M | Download |
| (MySQL-5.5.5_m3-1.sles11.src.rpm) | | MD5: e8a13c562e1846666367d211575f2f29 | |
| **Red Hat & Oracle Enterprise Linux 5 (Architecture Independent), RPM Package** | 5.5.5 | 21.0M | Download |
| (MySQL-5.5.5_m3-1.rhel5.src.rpm) | | MD5: dd01ca7e34be238d62eef7c1e3d4c3fc | |
| **SuSE Linux Enterprise Server 10 (Architecture Independent), RPM Package** | 5.5.5 | 21.0M | Download |
| (MySQL-5.5.5_m3-1.sles10.src.rpm) | | MD5: d5b5e453f40e4ecb35c2731f4ca9ab2b | |
| **Red Hat & Oracle Enterprise Linux 4 (Architecture Independent), RPM Package** | 5.5.5 | 21.0M | Download |
| (MySQL-5.5.5_m3-1.rhel4.src.rpm) | | MD5: 019c635b231625d206284e565079dccc | |
| **Generic Linux (Architecture Independent), Compressed TAR Archive** | 5.5.5 | 21.8M | Download |
| (mysql-5.5.5-m3.tar.gz) | | MD5: ad27f6561d6010c9346ffeca6de403fa | |

图 10-2 MySQL数据库图形管理工具下载

单击"Download"按钮后，进入的是下载镜像列表页面，由于mysql.com属美国网站，若在中国内地下载MySQL源码，访问美国的服务器会很慢，所以一般应选择亚洲的镜像站点，即选择"Asia"链接即可，如图10-3所示。



图 10-3 MySQL下载镜像列表

下载的源码包格式为tar.gz，它是一种压缩格式，在Linux用tar命令、Windows用WinRAR都可进行解压。解压缩后即可进行安装。MySQL源码包安装的操作步骤如下（如图10-4）

□□

- BUILD
- client
- cmake
- CMakeFiles
- cmd-line-utils
- config
- dbug
- Docs
- extra
- include
- libmysql
- libmysql_r
- libmysqld
- libservices
- man
- mysql-test
- mysys
- netware
- packaging
- plugin
- pstack
- regex
- scripts
- sql
- sql-bench
- sql-common
- storage
- strings
- support-files
- tests
- unittest
- vio
- win
- zlib

图 10-4 MySQL的插件式存储引擎

打开源代码压缩包，可以看到storage文件夹下面的文件夹，如图10-5所示。



图 10-5 存储引擎的文件夹名

从上面的文件夹名可以看出，目前的版本支持的存储引擎有以下几种：archive、blackhole、csv、fedorated、heap、ibmdb2i、myisam、innobase。
在MySQL 5.5版本以前，InnoDB Plugin作为默认的InnoDB存储引擎。打开MySQL 5.1安装目录，可以在安装目录下看到InnoDB插件，如图10-6所示。

图 10-6 MySQL 5.1的存储引擎目录

注意到这里有innobase和innodb_plugin两个目录，innobase目录保存的是InnoDB引擎的源代码，而innodb_plugin保存的是InnoDB Plugin引擎的源代码。如果要使用InnoDB Plugin，在编译MySQL时就应该指定使用innodb_plugin，而不是innobase。

[1] 见官网：http:‖www.mysql.com/downloads/mysql/。

# 10.2　InnoDB的逻辑存储

　　InnoDB存储引擎的逻辑存储结构和图10-7所示的其他数据库差不多，所有数据都被逻辑地存放在一个空间中，

- btr
- buf
- data
- dict
- dyn
- eval
- fil
- fsp
- fut
- ha
- handler
- ibuf
- include
- lock
- log
- mach
- mem
- mtr
- mysql-test
- os
- page
- pars
- que
- read
- rem
- row
- srv
- sync
- thr
- trx
- usr
- ut

表 10-7 InnoDB存储引擎源代码各模块及功能

□btr：B+树的实现。

□buf：缓冲池的实现，包括LRU算法、Flush刷新算法等。

□dict：InnoDB存储引擎中数据字典的实现。

□dyn：InnoDB存储引擎中动态数组的实现。

□fil：InnoDB存储引擎中文件的各种操作，包括文件的读取等。

□fsp：文件空间（file space）管理InnoDB存储引擎逻辑上对于文件的管理。

□ha：散列算法的实现。

□handler：实现对于MySQL的handler操作，即插件式存储引擎。

□ibuf：实现插入缓冲。

□include：InnoDB存储引擎的.h和.ic文件，存放所有数据结构的定义。

□lock：InnoDB存储引擎锁的实现，包括S锁、X锁及各种锁算法的实现。

□log：日志缓冲和重做日志文件的实现，包括写日志的算法及日志恢复的实现。

□mem：内存管理器的实现，用来分配内存和回收内存。

□mtr：小事务的实现。

□os：封装各个平台的系统调用。

□page：页的实现。

□row：对于各种类型的行操作实现。

□srv：对于InnoDB存储引擎参数的实现。

□sync：InnoDB存储引擎中对于Mutex量的实现。

□thr：InnoDB□□□□□□□□□□□□□□□□□

□trx：□□□□□□□□

□ut：□□□□□

# 10.3　MySQL 5.1□□□□□□□□InnoDB□□□□

## 10.3.1　Windows□□□□

□Windows□□□□□□□□□Visual Studion 2003□2005□2008□□□□□□MySQL□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□CMake□□□□□http://www.cmake.org□□□□

□bison□□□□□
http://gnuwin32.sourceforge.net/packages/bison.htm□□□□

□□□□□□□□□□configure.js□□□□□□□□□□□

```
C:\workdir□win\configure.js options
```

options□□□□□□□□□□□□

□WITH_INNOBASE_STORAGE_ENGINE□□□□InnoDB□□□□□□

□WITH_PARTITION_STORAGE_ENGINE□□□□□□□□

□WITH_ARCHIVE_STORAGE_ENGINE□□□□Archive□□□□□□

□WITH_BLACKHOLE_STORAGE_ENGINE□□□□Blackhole□□□□□□

□WITH_EXAMPLE_STORAGE_ENGINE□□□□Example□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□

□WITH_FEDERATED_STORAGE_ENGINE□□□□Federated□□□□□□

□WITH_NDBCLUSTER_STORAGE_ENGINE□□□□NDB Cluster□□□□□□

□□□□□□□□□□□□InnoDB□□□□□□□□□□□□□□□□□□□10-8□□□□

图　10-8　configure.js运行

上面的这个例子适合用于Visual Studio 2005及Visual Studio 2008。win目录下还有build-vsx.bat文件，表示Visual Studio的版本。build-vs8.bat表示Visual Studio 2005，build-vs8_x64.bat表示编译成64位的MySQL数据库，默认编译成32位的数据库。而对于Visual Studio 2008，可以看到如下过程所示的编译过程。

---

```
D:\Project\mysql-5.5.5-m3〉win\build-vs9.bat

--Check for working C compiler:C:/Program Files/Microsoft Visual Studio 9.0/VC/bin/cl.exe

--Check for working C compiler:C:/Program Files/Microsoft Visual Studio 9.0/VC/bin/cl.exe--works

--Detecting C compiler ABI info

--Detecting C compiler ABI info-done

--Check for working CXX compiler:C:/Program Files/Microsoft Visual Studio 9.0/VC/bin/cl.exe

--Check for working CXX compiler:C:/Program Files/Microsoft Visual Studio 9.0/VC/bin/cl.exe--works

--Detecting CXX compiler ABI info

--Detecting CXX compiler ABI info-done

--Check size of void*
```

--Check size of void*-done

SIZEOF_VOIDP=4

--Looking for include files HAVE_CXXABI_H

--Looking for include files HAVE_CXXABI_H-not found.

--Looking for include files HAVE_NDIR_H

--Looking for include files HAVE_NDIR_H-not found.

--Looking for include files HAVE_SYS_NDIR_H

--Looking for include files HAVE_SYS_NDIR_H-not found.

--Looking for include files HAVE_ASM_TERMBITS_H

--Looking for include files HAVE_ASM_TERMBITS_H-not found.

--Looking for include files HAVE_TERMBITS_H

--Looking for include files HAVE_TERMBITS_H-not found.

--Looking for include files HAVE_VIS_H

--Looking for include files HAVE_VIS_H-not found.

--Looking for include files HAVE_WCHAR_H

--Looking for include files HAVE_WCHAR_H-found

--Looking for include files HAVE_WCTYPE_H

--Looking for include files HAVE_WCTYPE_H-found

--Looking for include files HAVE_XFS_XFS_H

--Looking for include files HAVE_XFS_XFS_H-not found.

--Looking for include files CMAKE_HAVE_PTHREAD_H

--Looking for include files CMAKE_HAVE_PTHREAD_H-not found.

--Found Threads:TRUE

--Looking for pthread_rwlockattr_setkind_np

--Looking for pthread_rwlockattr_setkind_np-not found

--Performing Test HAVE_SOCKADDR_IN_SIN_LEN

--Performing Test HAVE_SOCKADDR_IN_SIN_LEN-Failed

--Performing Test HAVE_SOCKADDR_IN6_SIN6_LEN

--Performing Test HAVE_SOCKADDR_IN6_SIN6_LEN-Failed

--Cannot find wix 3,installer project will not be generated

--Configuring done

--Generating done

--Build files have been written to:D:/Project/mysql-5.5.5-m3

---

通过打开生成的MySQL.sln，接着通过编译工程，就可以对mysqld以及其他的组件进行调试，从而掌握MySQL的内部运行机制。

如果你已经迫不及待地想使用Visual Studio对数据库进行调试，可以参考第10-9节，了解InnoDB存储引擎master thread的运行方式。

图　10-9　查看master thread

# 10.3.2 Linux下的编译

Linux下同样可以使用Eclipse进行编译。在UNIX操作系统中（如Solaris、FreeBSD、MAC）都有相应的Eclipse版本，用户可以从http://www.eclipse.org/downloads/下载安装。Eclipse IDE for C/C++Developers版本已经对MySQL源码进行了配置，存放于/root/workspace/mysql-5.5.5-m3目录下。源码目录下有Make文件，Eclipse可以直接使用这个Make文件。

---

```
[root@xen-server mysql-5.5.5-m3]#BUILD/compile-amd64-debug-max-no-ndb-c
```

---

BUILD目录下有compile开头的一系列脚本，用户可以根据操作系统进行选择。64位的Linux，并且对源码进行Debug，因此选择compile-amd64-debug-max-no-ndb脚本，使用-c选项表示使用当前的Make文件进行编译。

接着启动Eclipse，选择新建C++工程，如图10-10所示。

图 10-10 应用程序C++代码

测试得到代码执行结果，并进入mysql_5_5_5数据库中查询数据，如图10-11所示。

图 10-11  生成的报表

单击Finish按钮完成报表的创建，最终结果如图10-12所示。

图 10-12 新建的C++项目

然后可以在左侧的Project Explorer窗口中的mysql_5_5_5项目树下看到所有的源码文件（/root/workspace/mysql-5.5.5-m3源码路径），如图10-13所示。

**New Folder**

**Folder**

Create a new folder resource.

Enter or select the parent folder:

mysql_5_5_5

mysql_5_5_5

Folder name: mysql-5.5.5-m3

<< Advanced

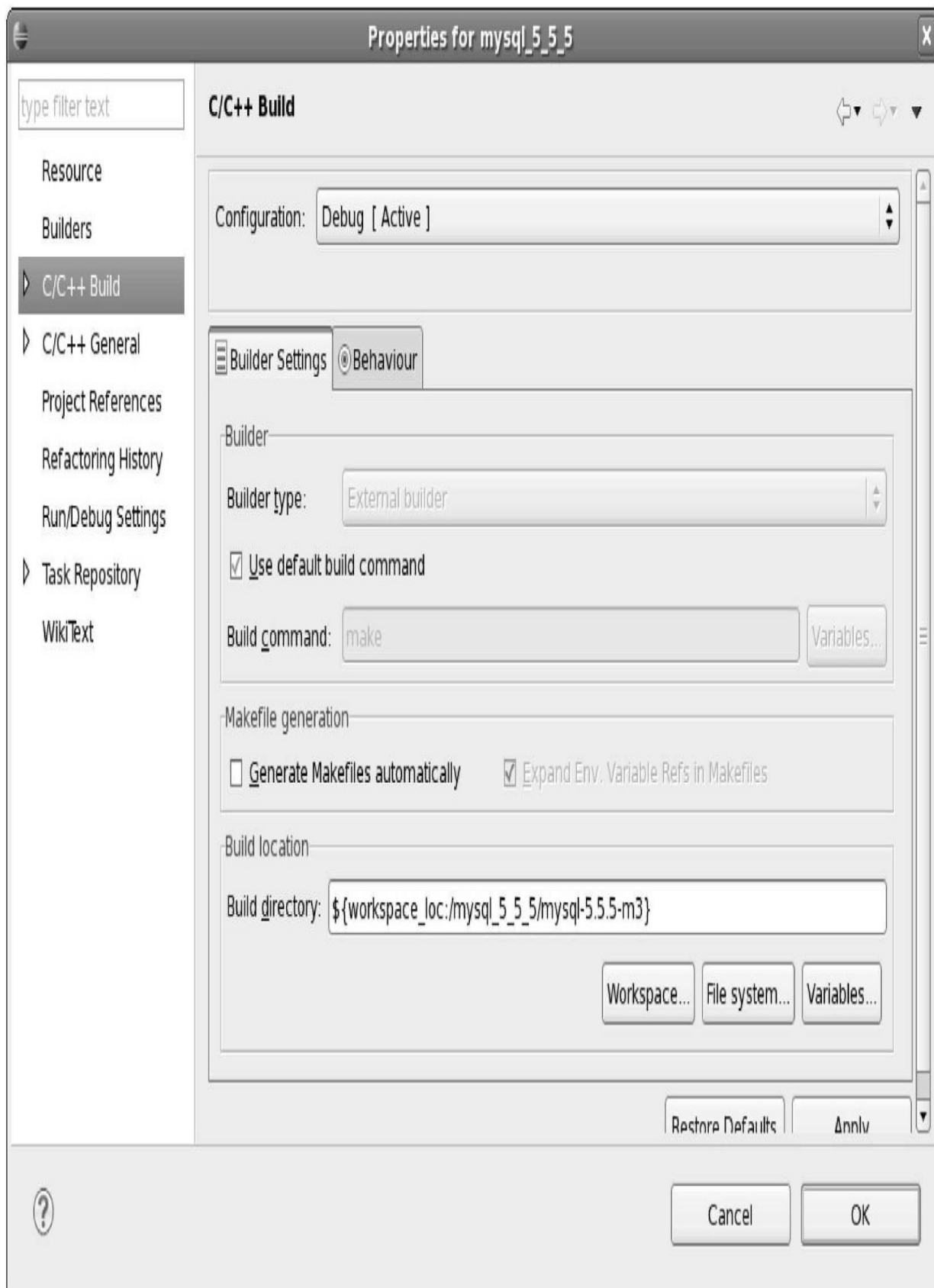☑ Link to folder in the file system

/root/workspace/mysql-5.5.5-m3     Browse...     Variables...

Cancel     Finish

图 10-13 选择组件

在该对话框中选择加入mysql_5_5_5库文件，然后单击C/C++Build菜单命令，在右侧
区域的Build directory文本框中选择编译的路径，如图10-14所示。

图 10-14 运动模糊

选择滤镜菜单下的模糊命令，其子菜单如图10-15所示。

图 10-15 属性设置

在选项卡中可以设置一些其他高级的服务器配置，但是一般不会直接设置mysqld的属性，这里不做介绍。单击"确定"按钮后，在Debug中执行启动命令，运行结果如图10-16所示。

图 10-16 Debug窗口

如果需要在运行程序时提供命令行参数，单击Arguments选项卡（图10-17），在

图 10-17  泵站布置

□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□□10-18所示。

图 10-18 　《Eclipse》操作界面

# 10.4　cmake环境编译源码InnoDB引擎调试

MySQL从版本5.5开始，编译源码从命令行cmake方式编译MySQL源码。如果想在 Mac OSX系统上想编译调试源码，这里介绍如何用xcode工具调MySQL源码的方式。

---

```
cd mysql-xxx

mkdir bld

cd bld

cmake..-G Xcode
```

---

编译完成后，MySQL的整个工程目录在bld目录下，然后再用cmake工具生成xcode工程文件，这样就可以用直接双击打开MySQL.xcodeproj工程，可以选择编译整个MySQL或者单独 InnoDB引擎，选择执行，如图10-19所示。

Running mysqld : ALL_BUILD

Project △ 257

Run　Stop　Scheme　Breakpoints　Editor　View　Organizer

ha_inn...　sync0s...　row0m...　trx0trx.h　ha_inn...　row0sel.c　btr0cur.c　sync0s...　srv0srv.c　+

MySQL ⟩ Sources ⟩ innobase ⟩ Source Files ⟩ srv0srv.c ⟩ srv_master_thread()

**mysqld**
Paused

▶ Thread 1
com.apple.main-thread

▼ Thread 2

　　0 srv_master_thread

　　1 _pthread_start

　　2 thread_start

▶ Thread 3

▶ Thread 4

▶ Thread 5

▶ Thread 6

▶ Thread 7

▶ Thread 8

▶ Thread 9

▶ Thread 10

▶ Thread 11

▶ Thread 12

▶ Thread 13

▶ Thread 14

▶ Thread 15

```c
2695    slot = srv_table_reserve_slot(SRV_MASTER);
2696
2697    srv_n_threads_active[SRV_MASTER]++;
2698
2699    mutex_exit(&kernel_mutex);
2700
2701 loop:
2702    /************************************************************/
2703    /* ---- When there is database activity by users, we cycle in this
2704    loop */
2705
2706    srv_main_thread_op_info = "reserving kernel mutex";
2707
2708    buf_get_total_stat(&buf_stat);
2709    n_ios_very_old = log_sys->n_log_ios + buf_stat.n_pages_read
2710        + buf_stat.n_pages_written;
2711    mutex_enter(&kernel_mutex);
2712
2713    /* Store the user activity counter at the start of this loop */
2714    old_activity_count = srv_activity_count;
2715
2716    mutex_exit(&kernel_mutex);
2717
2718    if (srv_force_recovery >= SRV_FORCE_NO_BACKGROUND) {      Thread 2: breakpoint 2.1
2719
2720        goto suspend_thread;
2721    }
2722
2723    /* ---- We run the following loop approximately once per second
2724    when there is database activity */
2725
2726    srv_last_log_flush_time = time(NULL);
2727
2728    /* Sleep for 1 second on entrying the for loop below the first time. */
2729    next_itr_time = ut_time_ms() + 1000;
2730
2731    for (i = 0; i < 10; i++) {
```

mysqld ⟩ Thread 2 ⟩ 0 srv_master_thread

Auto ⬩

▶ buf_stat (buf_pool_stat_t)

old_activity_count = (ulint) 0

srv_activity_count = (ulint) 0

srv_force_recovery = (ulint) 0

block->lock = Invalid Expression

图 10-19 Mac OS X中已经安装了的xcode、MySQL的开发

工具、cmake工具以及下载的MySQL源码包，接下来，使用cmake工具来编译和安装MySQL数据库软件。

## 10.5 小结

MySQL数据库的InnoDB存储引擎支持事务，事务处理是数据库应用开发中非常重要的内容。Visual Studio、Eclipse等开发工具都提供了对事务的支持，使得在应用程序中可以方便地使用事务来保证数据的一致性和完整性。此外，Oracle、Microsoft SQL Server、DB2等数据库管理系统也都支持事务处理。